



US011238239B2

(12) **United States Patent**
Liu et al.

(10) **Patent No.:** **US 11,238,239 B2**
(45) **Date of Patent:** **Feb. 1, 2022**

(54) **IN-CALL EXPERIENCE ENHANCEMENT FOR ASSISTANT SYSTEMS**

(71) Applicant: **Facebook Technologies, LLC**, Menlo Park, CA (US)

(72) Inventors: **Xiaohu Liu**, Bellevue, WA (US);
Baiyang Liu, Bellevue, WA (US);
Rajen Subba, San Carlos, CA (US)

(73) Assignee: **Facebook Technologies, LLC**, Menlo Park, CA (US)

(*) Notice: Subject to any disclaimer, the term of this patent is extended or adjusted under 35 U.S.C. 154(b) by 0 days.

(21) Appl. No.: **16/847,155**

(22) Filed: **Apr. 13, 2020**

(65) **Prior Publication Data**

US 2021/0120206 A1 Apr. 22, 2021

Related U.S. Application Data

(60) Provisional application No. 62/923,342, filed on Oct. 18, 2019.

(51) **Int. Cl.**

G06F 40/30 (2020.01)

G06F 9/54 (2006.01)

(Continued)

(52) **U.S. Cl.**

CPC **G06F 40/30** (2020.01); **G06F 3/011** (2013.01); **G06F 3/013** (2013.01); **G06F 9/485** (2013.01);

(Continued)

(58) **Field of Classification Search**

USPC 348/14.01

See application file for complete search history.

(56) **References Cited**

U.S. PATENT DOCUMENTS

6,462,660 B1 10/2002 Cannon et al.

7,124,123 B1 10/2006 Roskind

(Continued)

FOREIGN PATENT DOCUMENTS

AU 2017203668 1/2018

EP 2530870 12/2012

(Continued)

OTHER PUBLICATIONS

U.S. Appl. No. 14/593,723, filed Jan. 9, 2015, Colin Patrick Treseler.

(Continued)

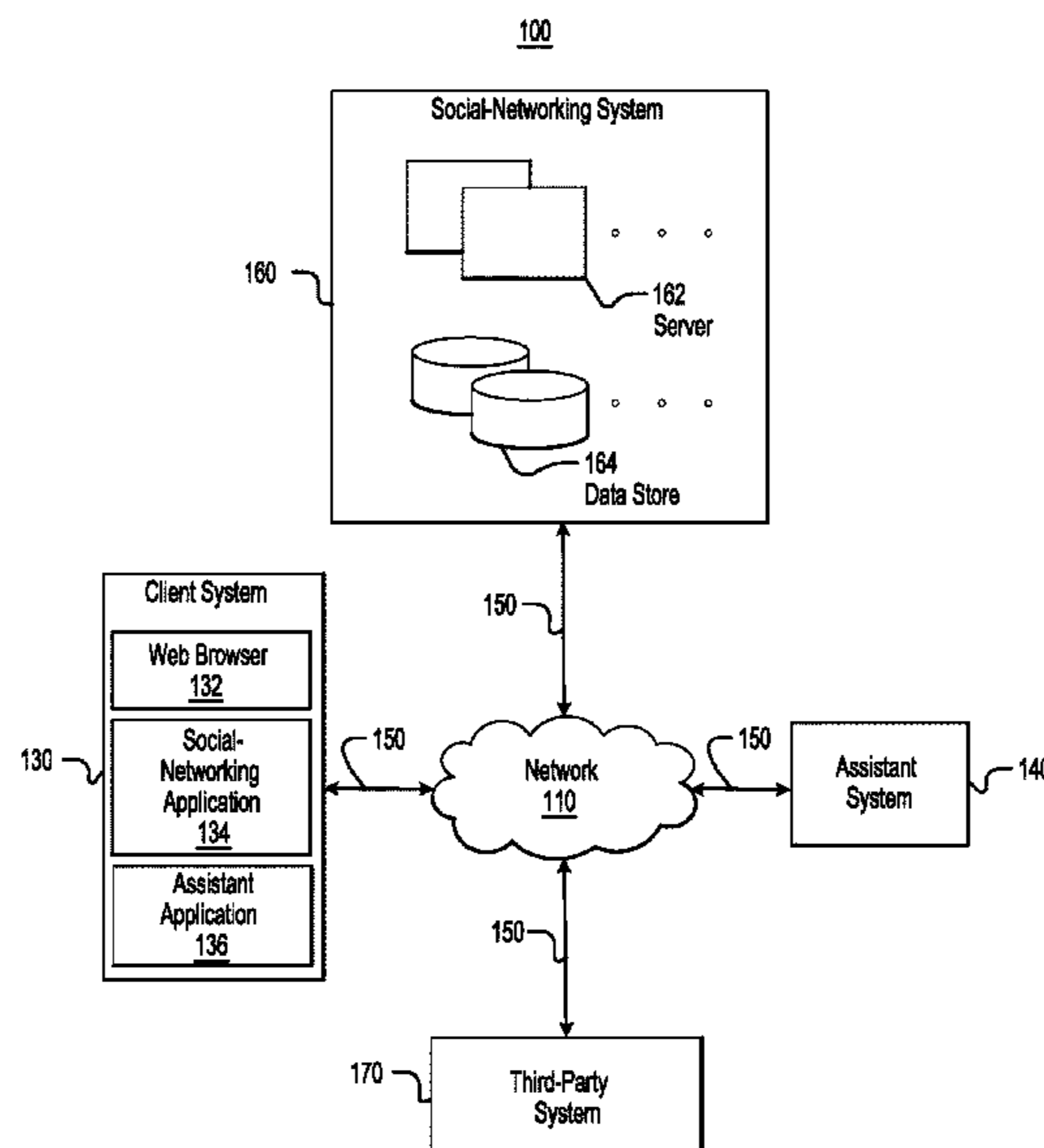
Primary Examiner — Maria El-Zoobi

(74) *Attorney, Agent, or Firm* — Baker Botts L.L.P.

(57) **ABSTRACT**

In one embodiment, a method includes establishing a video call between multiple client systems while persistently maintaining access to an assistant system during the video call. A request to be performed by the assistant system during the video call may then be received from a first client system; this request may reference one or more second users in the video call. An intent of the request and one or more user identifiers of these one or more second users referenced by the request may be determined, and the assistant system may be instructed to execute the request based on the determined intent and user identifiers. Finally, a response to the request may be sent to one more of the multiple client systems while maintaining the video call between these client systems.

20 Claims, 18 Drawing Sheets



(56)

References Cited

U.S. PATENT DOCUMENTS

2021/0064828 A1 3/2021 Johnson Premkumar et al.
 2021/0264913 A1 8/2021 Schramm et al.

FOREIGN PATENT DOCUMENTS

EP	3122001	1/2017
EP	3444814 A1	2/2019
WO	WO 2012/116241	8/2012
WO	2015179510 A1	11/2015
WO	2015183401 A1	12/2015
WO	WO 2016/195739	12/2016
WO	2017044163 A1	3/2017
WO	WO 2017/053208	3/2017
WO	WO 2017/116488	7/2017

OTHER PUBLICATIONS

- U.S. Appl. No. 15/808,638, filed Nov. 9, 2017, Ryan Brownhill.
 U.S. Appl. No. 15/949,011, filed Apr. 9, 2018, Jason Francis Harrison.
 U.S. Appl. No. 15/953,957, filed Apr. 16, 2018, Kemal El Moujahid.
 U.S. Appl. No. 15/966,455, filed Apr. 30, 2018, Scott Martin.
 U.S. Appl. No. 15/967,193, filed Apr. 30, 2018, Davide Testuggine.
 U.S. Appl. No. 15/967,239, filed Apr. 30, 2018, Vivek Natarajan.
 U.S. Appl. No. 15/967,279, filed Apr. 30, 2018, Fuchun Peng.
 U.S. Appl. No. 15/967,290, filed Apr. 30, 2018, Fuchun Peng.
 U.S. Appl. No. 15/967,342, filed Apr. 30, 2018, Vivek Natarajan.
 U.S. Appl. No. 16/011,062, filed Jun. 18, 2018, Jinsong Yu.
 U.S. Appl. No. 16/025,317, filed Jul. 2, 2018, Sonal Gupta.
 U.S. Appl. No. 16/036,827, filed Jul. 16, 2018, Emmanouil Koukoumidis.
 U.S. Appl. No. 16/038,120, filed Jul. 17, 2018, Jason Schissel.
 U.S. Appl. No. 16/048,049, filed Jul. 27, 2018, Markku Salkola.
 U.S. Appl. No. 16/048,072, filed Jul. 27, 2018, Markku Salkola.
 U.S. Appl. No. 16/048,101, filed Jul. 27, 2018, Markku Salkola.
 U.S. Appl. No. 16/053,600, filed Aug. 2, 2018, Vivek Natarajan.
 U.S. Appl. No. 16/057,414, filed Aug. 7, 2018, Jeremy Gillmor Kahn.
 U.S. Appl. No. 16/103,775, filed Aug. 14, 2018, Zheng Zhou.
 U.S. Appl. No. 16/107,601, filed Aug. 21, 2018, Rajesh Krishna Shenoy.
 U.S. Appl. No. 16/107,847, filed Aug. 21, 2018, Rajesh Krishna Shenoy.
 U.S. Appl. No. 16/118,169, filed Aug. 30, 2018, Baiyang Liu.
 U.S. Appl. No. 16/121,393, filed Sep. 4, 2018, Zheng Zhou.
 U.S. Appl. No. 16/127,173, filed Sep. 10, 2018, Zheng Zhou.
 U.S. Appl. No. 16/129,638, filed Sep. 12, 2018, Vivek Natarajan.
 U.S. Appl. No. 16/135,752, filed Sep. 19, 2018, Xiaohu Liu.
 U.S. Appl. No. 16/150,069, filed Oct. 2, 2018, Jiedan Zhu.
 U.S. Appl. No. 16/150,184, filed Oct. 2, 2018, Francislav P. Penov.
 U.S. Appl. No. 16/151,040, filed Oct. 3, 2018, Brian Nelson.
 U.S. Appl. No. 16/168,536, filed Oct. 23, 2018, Benoit F. Dumoulin.
 U.S. Appl. No. 16/176,0801, filed Oct. 31, 2018, Anusha Balakrishnan.
 U.S. Appl. No. 16/176,312, filed Oct. 31, 2018, Emmanouil Koukoumidis.
 U.S. Appl. No. 16/182,542, filed Nov. 6, 2018, Michael Robert Hanson.
 U.S. Appl. No. 16/183,650, filed Nov. 7, 2018, Xiaohu Liu.
 U.S. Appl. No. 16/192,538, filed Nov. 15, 2018, Emmanouil Koukoumidis.
 U.S. Appl. No. 16/222,923, filed Dec. 17, 2018, Jason Schissel.
 U.S. Appl. No. 16/222,957, filed Dec. 17, 2018, Emmanouil Koukoumidis.
 U.S. Appl. No. 16/229,828, filed Dec. 21, 2018, Xiaohu Liu.
 U.S. Appl. No. 16/247,439, filed Jan. 14, 2019, Xiaohu Liu.
 U.S. Appl. No. 16/264,173, filed Jan. 31, 2019, Ashwini Challa.
 U.S. Appl. No. 16/376,832, filed Apr. 5, 2019, Honglei Liu.
 U.S. Appl. No. 16/389,769, filed Apr. 19, 2019, Honglei Liu.
 U.S. Appl. No. 16/389,634, filed Apr. 19, 2019, Paul Anthony Crook.
 U.S. Appl. No. 16/389,738, filed Apr. 19, 2019, Fuchun Peng.
 U.S. Appl. No. 16/389,728, filed Apr. 19, 2019, William Crosby Present.
 U.S. Appl. No. 16/388,130, filed Apr. 18, 2019, Xiaohu Liu.
 U.S. Appl. No. 16/389,708, filed Apr. 19, 2019, William Crosby Present.
 U.S. Appl. No. 16/434,010, filed Jun. 6, 2019, Sergiu Dogaru.
 U.S. Appl. No. 16/552,559, filed Aug. 27, 2019, Seungwhan Moon.
 U.S. Appl. No. 16/557,055, filed Aug. 30, 2019, Seungwhan Moon.
 Mari Ostendorf, et al.: Human Language Technology: Opportunities and Challenges, IEEE International Conference on Acoustics, Speech, and Signal Processing, Mar. 23, 2005.
 U.S. Appl. No. 16/659,070, filed Oct. 21, 2019, Lisa Xiaoyi Huang.
 U.S. Appl. No. 16/659,203, filed Oct. 21, 2019, Lisa Xiaoyi Huang.
 U.S. Appl. No. 16/659,363, filed Oct. 21, 2019, Lisa Xiaoyi Huang.
 U.S. Appl. No. 16/659,419, filed Oct. 21, 2019, Lisa Xiaoyi Huang.
 U.S. Appl. No. 16/703,700, filed Dec. 4, 2019, Ahmed Aly.
 U.S. Appl. No. 16/733,044, filed Jan. 2, 2020, Francislav P. Penov.
 U.S. Appl. No. 16/741,630, filed Jan. 13, 2020, Paul Anthony Crook.
 U.S. Appl. No. 16/741,642, filed Jan. 13, 2020, Fuchun Peng.
 U.S. Appl. No. 16/742,769, filed Jan. 14, 2020, Xiaohu Liu.
 U.S. Appl. No. 16/742,668, filed Jan. 14, 2020, Xiaohu Liu.
 U.S. Appl. No. 16/790,497, filed Feb. 13, 2020, Yang Gao.
 U.S. Appl. No. 16/815,960, filed Mar. 11, 2020, Malik.
 U.S. Appl. No. 16/815,990, filed Mar. 11, 2020, Malik.
 U.S. Appl. No. 16/842,366, filed Apr. 7, 2020, Kamisetty.
 U.S. Appl. No. 62/660,876, filed Apr. 20, 2018, Anuj Kumar.
 U.S. Appl. No. 62/675,090, filed May 22, 2018, Michael Robert Hanson.
 U.S. Appl. No. 62/747,628, filed Oct. 18, 2018, Honglei Liu.
 U.S. Appl. No. 62/749,608, filed Oct. 23, 2018, Ashwini Challa.
 U.S. Appl. No. 62/750,746, filed Oct. 25, 2018, Honglei Liu.
 U.S. Appl. No. 62/923,342, filed Oct. 18, 2019, Michael Robert Hanson.
 Adam Poliak, et al.: Efficient, Compositional, Order-Sensitive n-gram Embeddings, Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: vol. 2, Short Papers, pp. 503-508, Valencia, Spain, Apr. 3-7, 2017.
 Chris Dyer, et al.: Recurrent Neural Network Grammars, Proceedings of NAACL-HLT 2016, pp. 199-209, San Diego, California, Jun. 12-17, 2016.
 Overview of Language Technology, <https://www.dfki.de/It/It-general.php>, Feb. 15, 2018.
 Jim Glass: A Brief Introduction to Automatic Speech Recognition, <http://www.cs.columbia.edu/~mcollins/6864/slides/asr.pdf>, Feb. 15, 2018.
 Speech Synthesis, https://en.wikipedia.org/wiki/Speech_synthesis, Feb. 15, 2018.
 Natural-language Understanding, https://en.wikipedia.org/wiki/Natural-language_understanding, Feb. 15, 2018.
 Conversational AI and the Road Ahead, <https://techcrunch.com/2017/02/25/conversational-ai-and-the-road-ahead/>, Feb. 15, 2018.
 What is Conversational AI? <https://blog.salemove.com/what-is-conversational-ai/>, Feb. 15, 2018.
 Question Answering, https://en.wikipedia.org/wiki/Question_answering, Feb. 15, 2018.
 Dialog Manager, https://en.wikipedia.org/wiki/Dialog_manager, Feb. 15, 2018.
 Dialog Management, <https://tutorials.botsfloor.com/dialog-management-799c20a39aad>, Feb. 15, 2018.
 Steeve Huang, Word2Vec and FastText Word Embedding with Gensim, <https://towardsdatascience.com/word-embedding-with-word2vec-and-fasttext-a209c1d3e12c>, Mar. 13, 2018.
 Ashwini Challa, et al.: Generate, Filter, and Rank: Grammaticality Classification for Production-Ready NLG Systems, arXiv: 1904.03279, Apr. 9, 2019.
 Honglei Liu, et al.: Explore-Exploit: A Framework for Interactive and Online Learning, arXiv: 1812.00116, Dec. 1, 2018.

(56)

References Cited

OTHER PUBLICATIONS

Turniški, Filip, et al. "Analysis of 3G and 4G download throughput in pedestrian zones." 2016 International Symposium ELMAR. IEEE, Sep. 12, 2016.

Dubin, Ran, et al. "Adaptation Logic for HTTP Dynamic Adaptive Streaming using Geo-Predictive Crowdsourcing," arXiv preprint arXiv: 1602.02030 (Feb. 5, 2016).

Chat Extensions, <https://developers.facebook.com/docs/messenger-platform/guides/chat-extensions>, Apr. 18, 2017.

Google Allo Makes Conversations Easier, Productive, and more Expressive, <https://www.trickyways.com/2016/05/google-allo-makes-conversations-easier-productive-expressive/>, May 19, 2016.

Agrawal, Aishwarya, et al. "VQA: Visual Question Answering." *International Journal of Computer Vision* 1.123 (2016): 4-31, Oct. 27, 2016.

Bast, Hannah, et al. "Easy access to the freebase dataset." Proceedings of the 23rd International Conference on World Wide Web. ACM, Apr. 7-11, 2014.

Bauer, Lisa, Yicheng Wang, and Mohit Bansal. "Commonsense for generative multi-hop question answering tasks." arXiv preprint arXiv: 1809.06309 (2018), Jun. 1, 2019.

Bordes, Antoine, et al. "Large-scale simple question answering with memory networks." arXiv preprint arXiv: 1506.02075 (Jun. 5, 2015).

Bordes, Antoine, et al. "Translating embeddings for modeling multi-relational data." *Advances in neural information processing systems*, Dec. 5, 2013.

Bordes, Antoine, Y-Lan Boureau, and Jason Weston. "Learning end-to-end goal-oriented dialog." arXiv preprint arXiv: 1605.07683 (2016), Mar. 30, 2017.

Choi, Eunsol, et al. "Quac: Question answering in context." arXiv preprint arXiv: 1808.07036 (Aug. 28, 2018).

Clark, Peter, et al. "Think you have solved question answering? try arc, the ai2 reasoning challenge." arXiv preprint arXiv: 1803.05457 (Mar. 14, 2018).

Conneau, Alexis, et al. "Supervised Learning of Universal Sentence Representations from Natural Language Inference Data." (Jul. 8, 2018).

Dalton, Jeffrey, Victor Ajayi, and Richard Main. "Vote Goat: Conversational Movie Recommendation," The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval. ACM, May 15, 2018.

de Vries, Harm, et al. "Talk the walk: Navigating newyork city through grounded dialogue." arXiv preprint arXiv: 1807.03367 (Dec. 23, 2018).

Dubey, Mohnish, et al. "Earl: Joint entity and relation linking for question answering over knowledge graphs." *International Semantic Web Conference*. Springer, Cham, Jun. 25, 2018.

Duchi, John, Elad Hazan, and Yoram Singer. "Adaptive subgradient methods for online learning and stochastic optimization." *Journal of Machine Learning Research* (Jul. 12, 2011): 2121-2159, Jul. 11, 2011.

Hudson, Drew A., and Christopher D. Manning. "Gqa: A new dataset for real-world visual reasoning and compositional question answering." *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, May 10, 2019.

Jiang, Lu, et al. "Memexqa: Visual memex question answering." arXiv preprint arXiv: 1708.01336 (Aug. 4, 2017).

Jung, Hyunwoo, et al. "Learning What to Remember: Long-term Episodic Memory Networks for Learning from Streaming Data." arXiv preprint arXiv: 1812.04227 (Dec. 11, 2018).

Kottur, Satwik, et al. "Visual coreference resolution in visual dialog using neural module networks." *Proceedings of the European Conference on Computer Vision (ECCV)*, Sep. 8-14, 2018.

Kumar, Ankit, et al. "Ask me anything: Dynamic memory networks for natural language processing." *International conference on machine learning*, Jan. 6, 2016.

Li, Jiwei, et al. "A persona-based neural conversation model." arXiv preprint arXiv: 1603.06155 (Jun. 8, 2016).

Pennington, Jeffrey, Richard Socher, and Christopher Manning. "Glove: Global vectors for word representation." *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, Oct. 25-29, 2014.

Rajpurkar, Pranav, et al. "Squad: 100,000+ questions for machine comprehension of text." arXivpreprint arXiv: 1606.05250 (Oct. 11, 2016).

Rajpurkar, Pranav, Robin Jia, and Percy Liang. "Know What You Don't Know: Unanswerable Questions for SQuAD." arXiv preprint arXiv: 1806.03822 (Jun. 11, 2018).

Reddy, Siva, Danqi Chen, and Christopher D. Manning. "Coqa: A conversational question answering challenge." *Transactions of the Association for Computational Linguistics* 7 (2019): 249-266, May 29, 2019.

Seo, Minjoon, et al. "Bidirectional attention flow for machine comprehension." arXiv preprint arXiv: 1611.01603 (Jun. 21, 2016).

Sukhbaatar, Sainbayar, Jason Weston, and Rob Fergus. "End-to-end memory networks." *Advances in neural information processing systems*, Nov. 24, 2015.

Tran, Ke, Arianna Bisazza, and Christof Monz. "Recurrent memory networks for language modeling." arXiv preprint arXiv: 1601.01272 (Apr. 22, 2016).

Wang, Peng, et al. "FVQA: Fact-based visual question answering." *IEEE transactions on pattern analysis and machine intelligence* 40.10 (2018): 2413-2427, Aug. 8, 2017.

Welbl, Johannes, Pontus Stenetorp, and Sebastian Riedel. "Constructing datasets for multi-hop reading comprehension across documents." *Transactions of the Association for Computational Linguistics* 6 (2018): 287-302, Jun. 11, 2018.

Salem, Yasser, Jun Hong, and Weiru Liu. "History-guided conversational recommendation." *Proceedings of the 23rd International Conference on World Wide Web*. ACM, Apr. 7-11, 2014.

Sun, Yueming, and Yi Zhang. "Conversational recommender system." *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval*. ACM, Jul. 8-12, 2018.

Wei, Wei, et al. "Airdialogue: An environment for goal-oriented dialogue research." *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, Nov. 4, 2018.

Weston, Jason, Sumit Chopra, and Antoine Bordes. "Memory networks." arXiv preprint arXiv: 1410.3916 (Nov. 29, 2014).

Williams, Jason D., Kavosh Asadi, and Geoffrey Zweig. "Hybrid code networks: practical and efficient end-to-end dialog control with supervised and reinforcement learning." arXiv preprint arXiv: 1702.03274 (Apr. 24, 2017).

Wu, Qi, et al. "Image Captioning and Visual Question Answering Based on Attributes and External Knowledge." arXiv preprint arXiv: 1603.02814 (Dec. 16, 2016).

Xu, Kun, et al. "Question answering on freebase via relation extraction and textual evidence." arXiv preprint arXiv: 1603.00957 (Jun. 9, 2016).

Yang, Zhilin, et al. "Hotpotqa: A dataset for diverse, explainable multi-hop question answering." arXiv preprint arXiv: 1809.09600 (Sep. 25, 2018).

Yin, Wenpeng, et al. "Simple question answering by attentive convolutional neural network." arXiv preprint arXiv: 1606.03391 (Oct. 11, 2016).

Zhang, Saizheng, et al. "Personalizing Dialogue Agents: I have a dog, do you have pets too?." arXiv preprint arXiv: 1801.07243 (Sep. 25, 2018).

Carlson, Andrew, et al. "Toward an architecture for never-ending language learning," *Twenty-Fourth AAAI Conference on Artificial Intelligence*, Jul. 5, 2010.

Chen, Chun-Yen, et al. "Gunrock: Building A Human-Like Social Bot By Leveraging Large Scale Real User Data." (2018).

Chen, Yun-Nung, William Yang Wang, and Alexander Rudnicky. "Jointly modeling inter-slot relations by random walk on knowledge graphs for unsupervised spoken language understanding." *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, May 31-Jun. 5, 2015.

Dettmers, Tim, et al. "Convolutional 2d knowledge graph embeddings." *Thirty-Second AAAI Conference on Artificial Intelligence*, Apr. 25, 2018.

(56)

References Cited

OTHER PUBLICATIONS

- Duchi, John, Elad Hazan, and Yoram Singer. "Adaptive subgradient methods for online learning and stochastic optimization." *Journal of Machine Learning Research* (Jul. 12, 2011): 2121-2159.
- Ganin, Yaroslav, et al. "Domain-adversarial training of neural networks." *The Journal of Machine Learning Research* 17.1 (2016): 2096-2030, Jan. 1, 2016.
- Ghazvininejad, Marjan, et al. "A knowledge-grounded neural conversation model." *Thirty-Second AAAI Conference on Artificial Intelligence*, Apr. 27, 2018.
- He, He, et al. "Learning symmetric collaborative dialogue agents with dynamic knowledge graph embeddings." *arXiv preprint arXiv: 1704.07130* (Apr. 14, 2017).
- Henderson, Matthew, Blaise Thomson, and Jason D. Williams. "The second dialog state tracking challenge." *Proceedings of the 15th Annual Meeting of the Special Interest Group on Discourse and Dialogue (SIGDIAL)*, Jun. 18-20, 2014.
- Kartsaklis, Dimitri, Mohammad Taher Pilehvar, and Nigel Collier. "Mapping text to knowledge graph entities using multi-sense Istm." *arXiv preprint arXiv: 1808.07724* (Aug. 23, 2018).
- Lao, Ni, Tom Mitchell, and William W. Cohen. "Random walk inference and learning in a large scale knowledge base." *Proceedings of the Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Jul. 27, 2011.
- Long, Yinong, et al. "A knowledge enhanced generative conversational service agent." *DSTC6 Workshop*, Dec. 2017.
- Miller, Alexander H., et al. "Parlai: A dialog research software platform." *arXiv preprint arXiv: 1705.06476* (May 18, 2017).
- Moon, Seungwhan, and Jaime G. Carbonell. "Completely Heterogeneous Transfer Learning with Attention-What And What Not To Transfer." *IJCAI*, Aug. 19, 2017.
- Moon, Seungwhan, Suyoun Kim, and Haohan Wang. "Multimodal transfer deep learning with applications in audio-visual recognition." *arXiv preprint arXiv: 1412.3121* (Dec. 9, 2014).
- Moon, Seungwhan, Leonardo Neves, and Vitor Carvalho. "Multimodal named entity recognition for short social media posts." *arXiv preprint arXiv: 1802.07862* (Feb. 22, 2018).
- Moon, Seungwhan, Leonardo Neves, and Vitor Carvalho. "Zeroshot Multimodal Named Entity Disambiguation for Noisy Social Media Posts." *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (vol. 1: Long Papers)*, Jul. 15-20, 2018.
- Nickel, Maximilian, Lorenzo Rosasco, and Tomaso Poggio. "Holographic embeddings of knowledge graphs." *Thirtieth Aaai conference on artificial intelligence*, Mar. 2, 2016.
- Parthasarathi, Prasanna, and Joelle Pineau. "Extending neural generative conversational model using external knowledge sources." *arXiv preprint arXiv: 1809.05524* (Sep. 14, 2018).
- Salem, Yasser, Jun Hong, and Weiru Liu. "History-guided conversational recommendation." *Proceedings of the 23rd International Conference on World Wide Web*. ACM, Apr. 7, 2014.
- Shah, Pararth, et al. "Bootstrapping a neural conversational agent with dialogue self-play, crowdsourcing and on-line reinforcement learning." *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, vol. 3 (Industry Papers)*, Jun. 1-6, 2018.
- Sutskever, Ilya, Oriol Vinyals, and Quoc V. Le. "Sequence to sequence learning with neural networks." *Advances in neural information processing systems*, Sep. 10, 2014.
- Wang, Zhen, et al. "Knowledge graph embedding by translating on hyperplanes." *Twenty-Eighth AAAI conference on artificial intelligence*, Jun. 21, 2014.
- Yang, Zichao, et al. "Hierarchical attention networks for document classification." *Proceedings of the 2016 conference of the North American chapter of the association for computational linguistics: human language technologies*, Jun. 12-17, 2016.
- Young, Tom, et al. "Augmenting end-to-end dialogue systems with commonsense knowledge." *Thirty-Second AAAI Conference on Artificial Intelligence*, Apr. 26, 2018.
- Zhang, Saizheng, et al. "Personalizing Dialogue Agents: I have a dog, do you have pets too?." *arXiv preprint arXiv: 1801.07243* (Jan. 22, 2018).
- Bui, Duc, et al., "Federated User Representation Learning" *arXiv: 1909.12535*, Sep. 27, 2019.
- Goetz, Jack, et al. "Active Federated Learning" *arXiv: 1909.12641*, Sep. 27, 2019.
- Hsiao, Wei-Lin, et al. "Fashion++: Minimal Edits for Outfit Improvement." *arXiv preprint arXiv: 1904.09261* (Apr. 19, 2019).
- Dinan, Emily, et al. "Advances in Conversational AI" https://ai.facebook.com/blog/advances-in-conversational-ai/?_xts=%5b0%5d=68_ARDgZpslcbW2Y4dGWBFBIBBfrsZkeNMXeTFXLveffyaOCRJ0INA80NQAj9Y6urka2DI6EQcbA0JoTxUuSGUFT-BkfYahB6ILnX-UMQR5WBIRXKYbJ43SS6THchGawfWiM4ESBHe_Ok7V7IUT97zwwqpCz26VR0EQTuELalEkfW1sb7BGN16RGOMeIOCRC38TiqG3U-3Vk0Mns4L-esrTNWVC3RoyMpYGOFKTwMWb2q8yHbkVod9ZwDKi6XC01CIbVIwa_BAz3zINQR-FV4z-lkOf7M-xGMuXDbTjDgf7nhCeQmY4Afa-Dcf6OVbG2dKgFXztn3Jqi0iVPBw_nVko&_in_=UK-R, Aug. 2, 2019.
- Ott, Myle, et al. "New advances in natural language processing to better connect people" https://ai.facebook.com/blog/new-advances-in-natural-language-processing-to-better-connect-people/?_xts=%5b0%5d=68_ARBpsX-0s8sV0sN3kxbWpoyzVrkSOpof5SCANCyVPB6BtolxwZPobEfGIXdGEOofVPVTA3-LJPx6LICOHs5_Kqixd4ZXIiEssjj04COGloA0SmwZeEDo2tSV4hCmYGITKMotQzPWB6QA9iS_c5_13tm4jANzArPOC6M9tzzXxfmixtWiYv-Zkvc2dJ-9MFVvXrv2vxijyqLOat3oKUvmwPO5Gzny2yiqiE6JT8uXmXvBJgXSJaxHZYvhoaP6sEWsMt6LOBtC_DzR7nIibeSBAIdbh2EYVHRLJzdp7flnuepGacg7ZmKHJ4P235Qyi6fy06TFynBX7x8AeRGAqI5dbkg&_tn_=UK-R, Aug. 14, 2019.
- Oord, Aaron van den, et al. "Wavenet: A generative model for raw audio." *arXiv preprint arXiv: 1609.03499* (Sep. 12, 2016).
- Bahdanau, Dzmitry, Kyunghyun Cho, and Yoshua Bengio. "Neural machine translation by jointly learning to align and translate." *arXiv preprint arXiv: 1409.0473* (Sep. 1, 2014).
- Vaswani, Ashish, et al. "Attention is all you need." *Advances in neural information processing systems*. 2017.
- Shen, Jonathan, et al. "Natural tts synthesis by conditioning wavenet on mel spectrogram predictions." *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, Apr. 15, 2018.
- Wang, Yuxuan, et al. "Uncovering latent style factors for expressive speech synthesis." *arXiv preprint arXiv: 1711.00520* (Nov. 1, 2017).
- Skerry-Ryan, R. J., et al. "Towards end-to-end prosody transfer for expressive speech synthesis with tacotron." *arXiv preprint arXiv: 1803.09047* (Mar. 24, 2018).
- Wang, Yuxuan, et al. "Style tokens: Unsupervised style modeling, control and transfer in end-to-end speech synthesis." *arXiv preprint arXiv: 1803.09017* (Mar. 23, 2018).
- Locatello, Francesco, et al. "Challenging common assumptions in the unsupervised learning of disentangled representations." *arXiv preprint arXiv: 1811.12359* (Nov. 29, 2018).
- Tits, Noé, et al. "Visualization and interpretation of latent spaces for controlling expressive speech synthesis through audio analysis." *arXiv preprint arXiv: 1903.11570* (Mar. 27, 2019).
- Renter, Gustav Eje, et al. "Principles for Learning Controllable TTS from Annotated and Latent Variation." *INTERSPEECH*. 2017.
- Hodari, Zack, et al. "Learning Interpretable Control Dimensions for Speech Synthesis by Using External Data." *Interspeech*. 2018.
- Busso, Carlos, et al. "IEMOCAP: Interactive emotional dyadic motion capture database." *Language resources and evaluation* 42.4 (Dec. 1, 2008): 335.
- Simon King, Lovisa Wihlborg, and Wei Guo, "The blizzard challenge 2017," in *Proc. Blizzard Challenge*, 2017, vol. 2017, pp. 1-17.
- Yoon, Seunghyun, Seokhyun Byun, and Kyommm Jung. "Multimodal speech emotion recognition using audio and text." *2018 IEEE Spoken Language Technology Workshop (SLT)*. IEEE, Dec. 18, 2018.

(56)

References Cited

OTHER PUBLICATIONS

Banse, Rainer, and Klaus R. Scherer. "Acoustic profiles in vocal emotion expression." *Journal of personality and social psychology* 70.3 (1996): 614, 1996.

Yamagishi, Junichi, et al. "Modeling of various speaking styles and emotions for HMM-based speech synthesis." *Eighth European Conference on Speech Communication and Technology*. 2003.

Tachibana, Makoto, et al. "HMM-based speech synthesis with various speaking styles using model interpolation." *Speech Prosody 2004, International Conference*. 2004.

Yamagishi, Junichi, et al. "Speaking style adaptation using context clustering decision tree for HMM-based speech synthesis." *2004 IEEE International Conference on Acoustics, Speech, and Signal Processing*. vol. 1. IEEE, May 17, 2004.

Lee, Chi-Chun, et al. "Emotion recognition using a hierarchical binary decision tree approach." *Speech Communication* 53.9-10 (2011): 1162-1171, Nov. 1, 2011.

Mower, Emily, Maja J. Mataric, and Shrikanth Narayanan. "A framework for automatic human emotion classification using emotion profiles." *IEEE Transactions on Audio, Speech, and Language Processing* 19.5 (2010): 1057-1070, Sep. 27, 2010.

Han, Kun, Dong Yu, and Ivan Tashev. "Speech emotion recognition using deep neural network and extreme learning machine." *Fifteenth annual conference of the international speech communication association*. 2014.

Trigeorgis, George, et al. "Adieu features? end-to-end speech emotion recognition using a deep convolutional recurrent network." *2016 IEEE international conference on acoustics, speech and signal processing (ICASSP)*. IEEE, Mar. 20, 2016.

Li, Yanghao, et al. "Adaptive batch normalization for practical domain adaptation," *Pattern Recognition* 80 (2018): 109-117, Aug. 1, 2018.

Scherer, Klaus R., et al. "Vocal cues in emotion encoding and decoding." *Motivation and emotion* 15.2 (1991): 123-148, Jun. 1, 1991.

Yang Gao, "Demo for 'interactive text-to-speech via semi-supervised style transfer learning'," <https://github.com/Yolanda-Gao/Interactive-Style-TTS>, 2019, Accessed: Oct. 21, 2019.

Anonymous, "Make it Famous (Experiential photography platform)," Jun. 1, 2019 [Retrieved on Nov. 23, 2020], pp. 1-9, XP055753168, Retrieved from the Internet: URL: https://www.miracam.com/wp-content/themes/mira/assets/pdf/Mira_SpecSheet_201906.pdf.

Betermieux S., et al., "Finalizing Dialog Models at Runtime," Big Data Analytics in the Social and Ubiquitous Context, Jul. 16, 2007, 15 Pages, XP047417162.

Billsus D., et al., "Improving Proactive Information Systems," 2005 International Conference on Intelligent User Interfaces, IUI 05, Jan. 9-12, 2005, pp. 159-166, XP058319582.

Bonnington C., "Google Clips Smart Camera isn't Smart Enough, But its Aims are Still Worth Considering," *Slate* [Online], Feb. 27, 2018 [Retrieved on Nov. 23, 2020], pp. 1-3, XP055753206, Retrieved from the Internet: URL <https://slate.com/technology/2018/02/google-clips-smart-camera-isnt-smart-enough-but-its-aims-are-still-worth-considering.html>.

Constantinides P.C., et al., "A Schema Based Approach to Dialog Control," 5th International Conference on Spoken Language Processing, Oct. 1, 1998, 4 Pages, XP007000460.

Co-pending U.S. Appl. No. 16/153,574, inventors Jason; Harrison et al., filed on Oct. 5, 2018.

Co-Pending U.S. Appl. No. 16/914,966, inventor Noam; Behar, filed on Jun. 29, 2020.

Co-Pending U.S. Appl. No. 16/917,664, inventors Xiaohu; Liu et al., filed on Jun. 30, 2020.

Co-Pending U.S. Appl. No. 16/921,665, inventors Honglei; Liu et al., filed on Jul. 6, 2020.

Co-Pending U.S. Appl. No. 16/998,423, inventors Armen; Aghajanyan et al., filed on Aug. 20, 2020.

Co-Pending U.S. Appl. No. 17/006,260, inventors William; Presant et al., filed on Aug. 28, 2020.

Co-Pending U.S. Appl. No. 17/006,339, inventors Shivani; Poddar et al., filed on Aug. 28, 2020.

Co-Pending U.S. Appl. No. 17/006,377, inventors Shivani; Poddar et al., filed on Aug. 28, 2020.

Co-Pending U.S. Appl. No. 17/009,542, inventor Satwik; Kottur, filed on Sep. 1, 2020.

Co-Pending U.S. Appl. No. 17/035,253, inventors Piyush; Khemka et al., filed on Sep. 28, 2020.

Co-Pending U.S. Appl. No. 17/120,013, inventors Fadi; Botros et al., filed on Dec. 11, 2020.

Co-Pending U.S. Appl. No. 17/136,636, inventors Michael; Greenberg et al., filed on Dec. 29, 2020.

Co-pending U.S. Appl. No. 17/139,363, inventors Daniel; Cheng et al., filed on Dec. 31, 2020.

Co-pending U.S. Appl. No. 17/186,459, inventors Bing; Liu et al., filed on Feb. 26, 2021.

Daha F.Z., et al., "Deep Neural Architecture with Character Embedding for Semantic Frame Detection," 2019 IEEE 13th International Conference on Semantic Computing (ICSC), Jan. 30, 2019, pp. 302-307, XP033529242.

Dauphin Y.N., et al., "Language Modeling with Gated Convolutional Networks," Cornell University Library, NY 14853, Dec. 23, 2016, 9 Pages, XP080742751.

Duong L., et al., "An Adaptable Task-Oriented Dialog System for Stand-alone Embedded Devices," Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, System Demonstrations, Jul. 28, 2019, pp. 19-57, XP055767185.

Elgan M., "How Lifelogging Will Become Easy and Automatic," *Computerworld* [Online], Nov. 19, 2016 [Retrieved on Nov. 23, 2020], pp. 1-7, XP055753230, Retrieved from the Internet: URL: <https://www.computerworld.com/article/3143115/how-lifelogging-will-become-easy-and-automatic.html>.

Fiscus J.G., et al., "A Post-processing System to Yield Reduced Word Error Rates: Recognizer Output Voting Error Reduction (ROVER)," *Automatic Speech Recognition and Understanding*, Dec. 14-17, 1997, pp. 347-354, XP010267529.

Hubara I., et al., "Quantized Neural Networks: Training Neural Networks with Low Precision Weights and Activations," *Journal of Machine Learning Research* [Online], Sep. 22, 2016 [Retrieved on Jan. 20, 2021], vol. 18, 30 Pages, XP055611103, Retrieved from the Internet: URL: <https://arxiv.org/pdf/1609.07061.pdf>.

International Search Report and Written Opinion for International Application No. PCT/US2020/052866, dated Dec. 9, 2020, 11 Pages.

International Search Report and Written Opinion for International Application No. PCT/US2020/052955, dated Nov. 27, 2020, 12 Pages.

International Search Report and Written Opinion for International Application No. PCT/US2020/052966, dated Dec. 14, 2020, 12 Pages.

International Search Report and Written Opinion for International Application No. PCT/US2020/052969, dated Dec. 9, 2020, 11 Pages.

International Search Report and Written Opinion for International Application No. PCT/US2020/054604, dated Jan. 29, 2021, 17 Pages.

International Search Report and Written Opinion for International Application No. PCT/US2020/056144, dated Mar. 30, 2021, 8 Pages.

International Search Report and Written Opinion for International Application No. PCT/US2020/056188, dated Dec. 9, 2020, 12 Pages.

Kim Y., et al., "Character-Aware Neural Language Models," Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence [Online], Oct. 16, 2015 [Retrieved on Apr. 12, 2018], 9 Pages, XP055466626, Retrieved from the Internet URL: <https://arxiv.org/pdf/1508.06615.pdf>.

Laddha A., et al., "Understanding Chat Messages for Sticker Recommendation in Messaging Apps," Cornell University Library, NY 14853, Feb. 7, 2019, 8 Pages, XP081537983.

Mamou J., et al., "System Combination and Score Normalization for Spoken Term Detection," IEEE International Conference on

(56)

References Cited

OTHER PUBLICATIONS

Acoustics and Speech and Signal Processing, Proceedings 1999, May 26, 2013, pp. 8272-8276, XP032508928.

Martin S., et al., "MuDoCo: Corpus for Multidomain Coreference Resolution and Referring Expression Generation," In Proceedings of the 12th Conference on Language Resources and Evaluation Conference, May 16, 2020, pp. 104-111.

Moon S., et al., "Situated and Interactive Multimodal Conversations," Facebook AI Research, Jun. 2, 2020, pp. 1-16.

"Social Context Reminder", An IP.com Prior Art Database Technical Disclosure [Online], Aug. 4, 2017 [Retrieved on Oct. 27, 2020], 3 Pages, Retrieved from Internet: URL: <https://priorart.ip.com/IPCOM/000250582>.

Ostendorf M., et al., "Continuous-Space Language Processing: Beyond Word Embeddings," Springer International Publishing, Sep. 21, 2016, 13 Pages, XP047356965.

Pavel M., et al., "Behavioral Informatics and Computational Modeling in Support of Proactive Health Management and Care," IEEE Transactions on Biomedical Engineering, IEEE Service Center, Piscataway, NJ, USA, Dec. 1, 2015, vol. 62 (12), 32 Pages, XP011590079.

Planells J., et al., "A Multi-domain Dialog System to Integrate Heterogeneous Spoken Dialog Systems," Interspeech [Online], Aug. 25, 2013, 5 Pages, XP055650758, Retrieved from the Internet: URL: [https://www.isca-speech.org/archive/archive_papers/interspeech_2013/i13_1891 .pdf](https://www.isca-speech.org/archive/archive_papers/interspeech_2013/i13_1891.pdf).

Sarikaya R., "The Technology Behind Personal Digital Assistants: An Overview of the System Architecture and Key Components," IEEE Signal Processing Magazine, Jan. 1, 2017, vol. 34 (1), pp. 67-81, XP011639190.

Shiee N., et al., "Take Your Best Selfie Automatically, with Photobooth on Pixel 3," Google AI Blog [Online], Apr. 1, 2019 [Retrieved on 2020-11-24], pp. 1-4, XP055753559, Retrieved from the Internet: URL: <https://ai.googleblog.com/2019/04/take-your-best-selfie-automatically.html>.

Tepper N., et al., "Collabot: Personalized Group Chat Summarization," In Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining, Feb. 5, 2018, pp. 771-774.

Nang Y., et al., "Dialogue Intent Classification with Character-CNN-BGRU Networks," Multimedia Tools and Applications, Jun. 11, 2019, vol. 79 (8), pp. 4553-4572, XP037048798.

Ku P., et al., "Convolutional Neural Network Based Triangular CRF for Joint Intent Detection and Slot Filling," 2013 IEEE Workshop on Automatic Speech Recognition and Understanding, Dec. 8, 2013, pp. 78-83, XP032544418.

Yeung K.F., et al., "A Proactive Personalised Mobile Recommendation System Using Analytic Hierarchy Process and Bayesian Network," Journal of Internet Services and Applications, Jul. 20, 2012, vol. 3 (2), pp. 195-214, XP055754649.

Zhu M.H., et al., "To Prune, or Not to Prune: Exploring the Efficacy of Pruning for Model Compression," Cornell University Library, NY 14853, Oct. 5, 2017, 11 Pages, XP081283371.

* cited by examiner

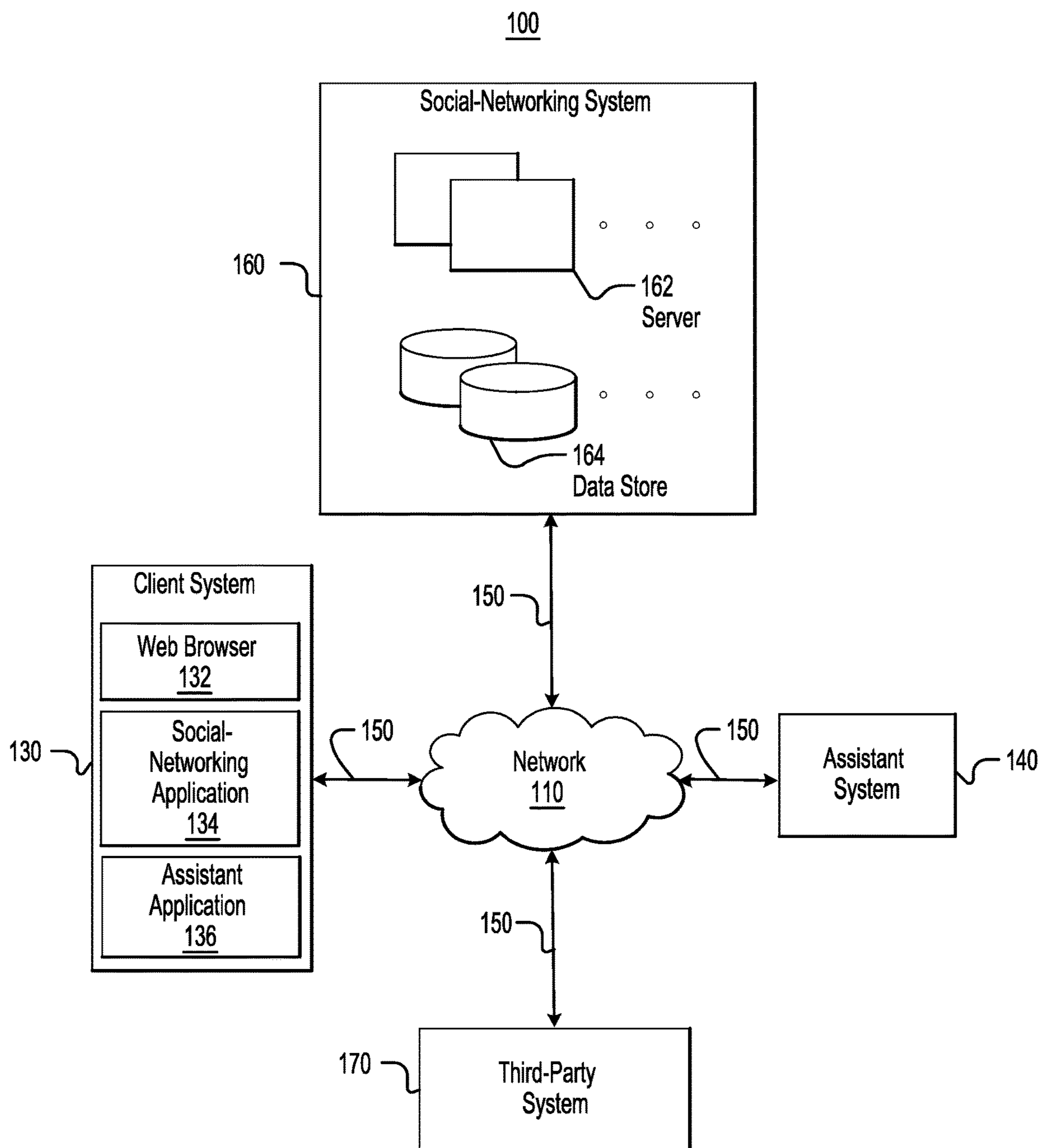


FIG. 1

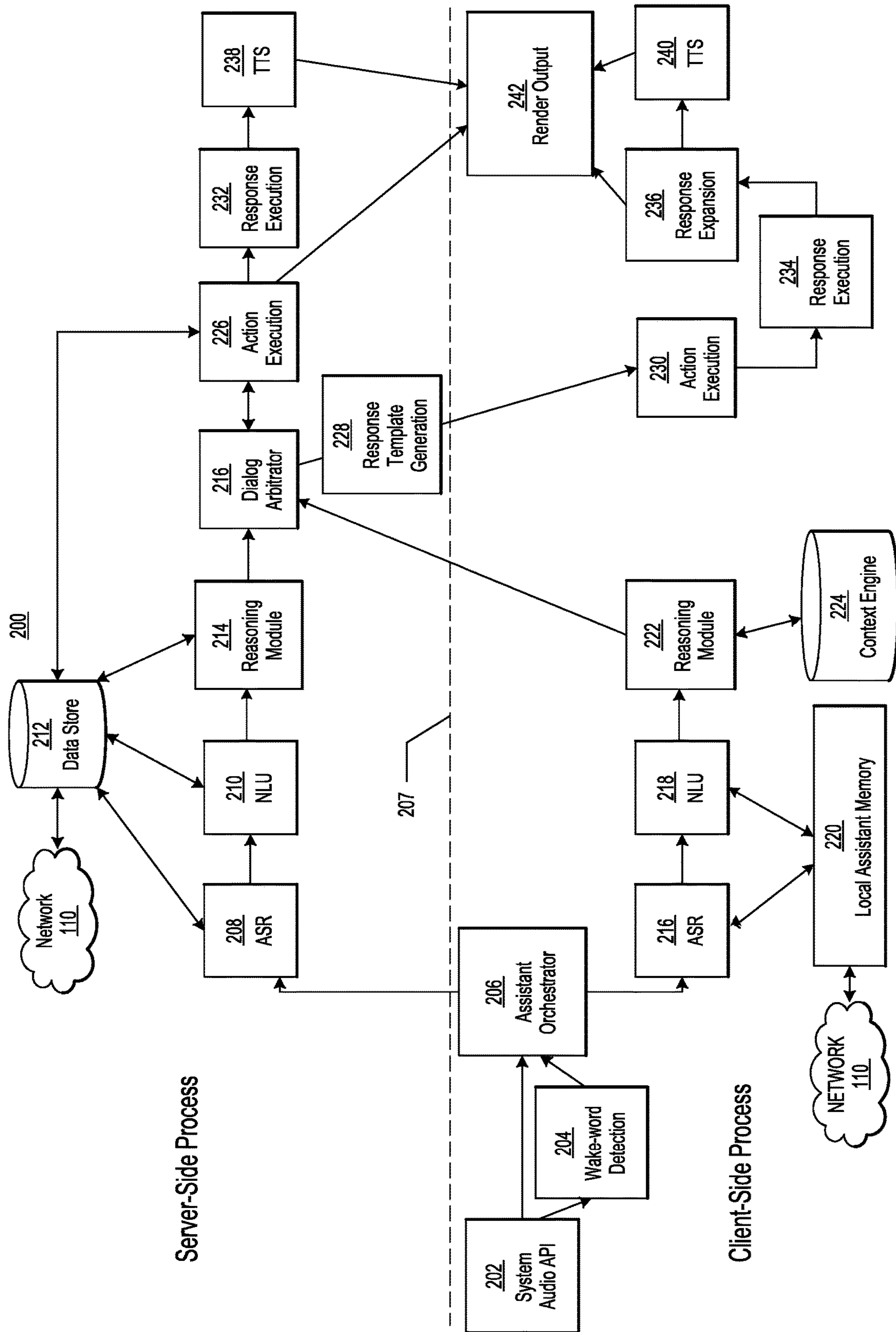


FIG. 2

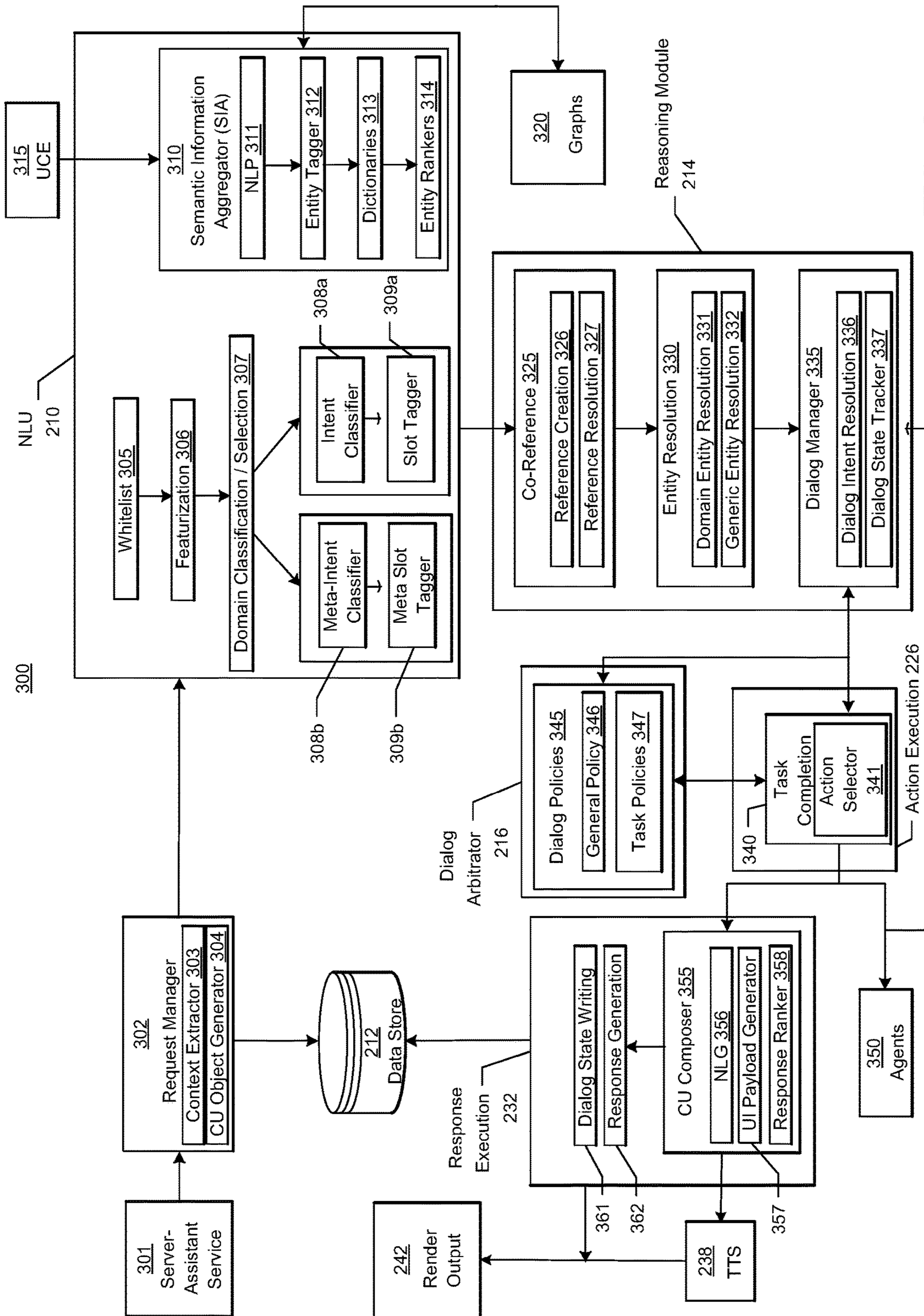


FIG. 3

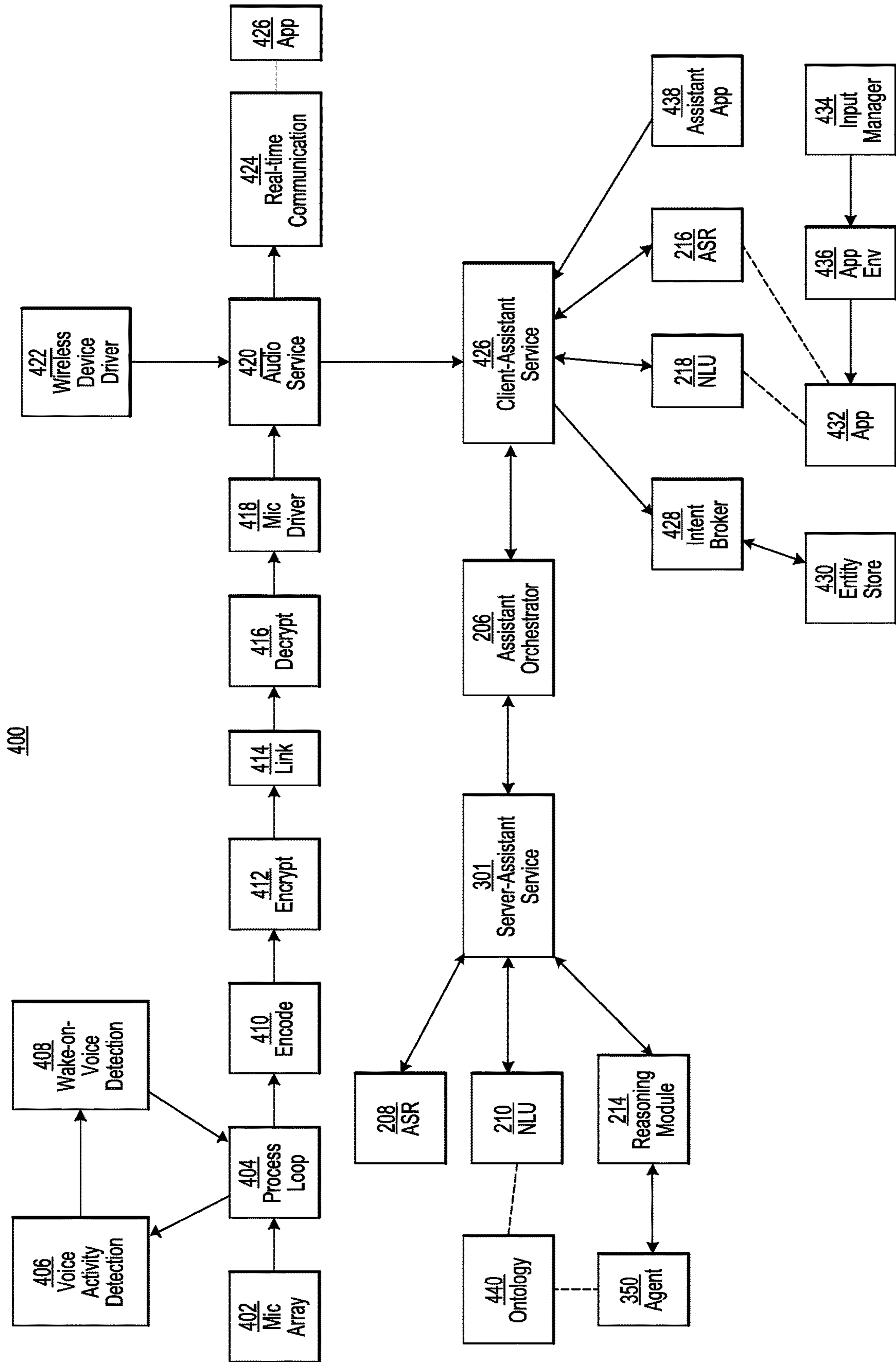


FIG. 4

Multimodal Architecture 500

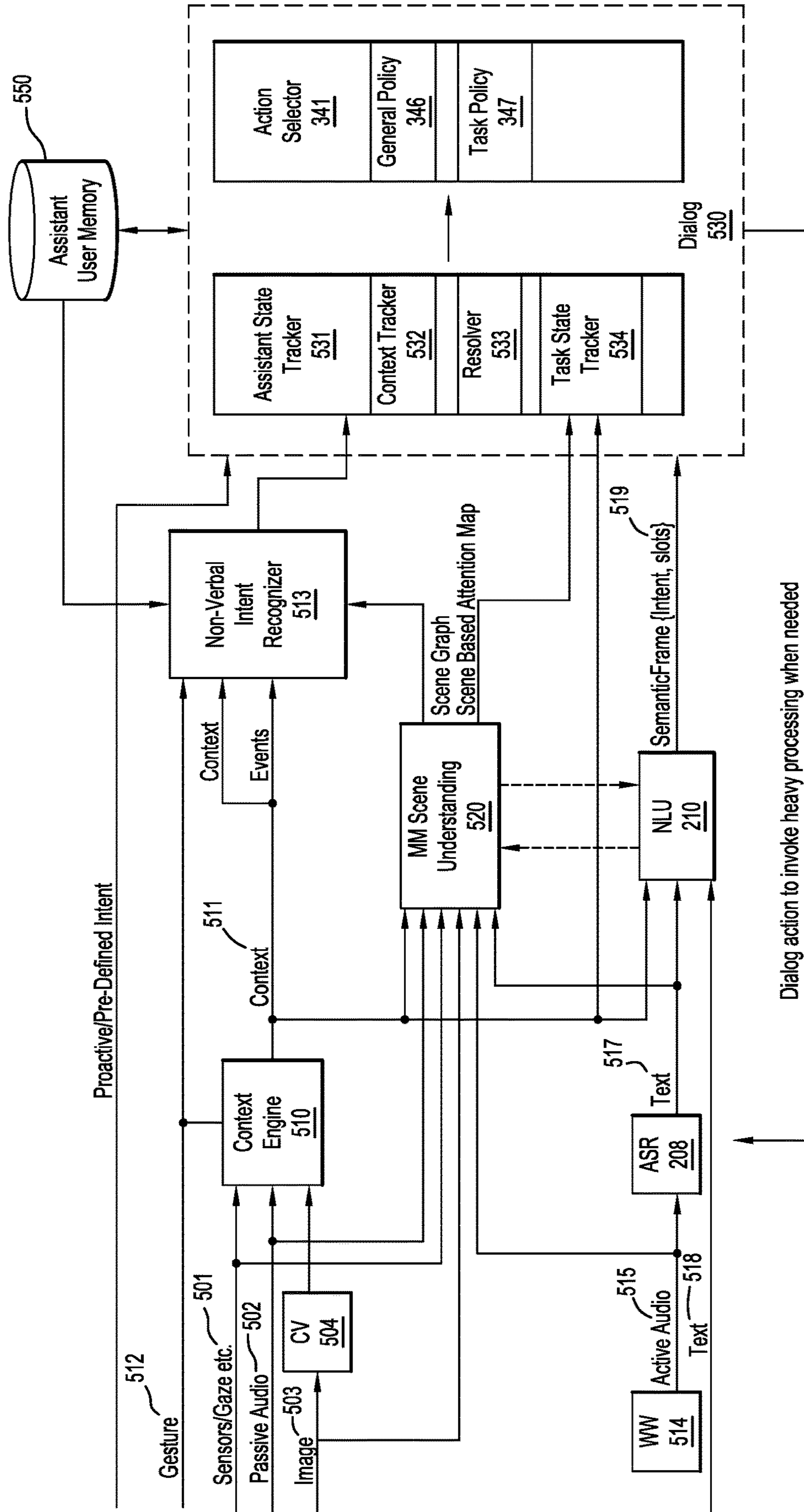


FIG. 5

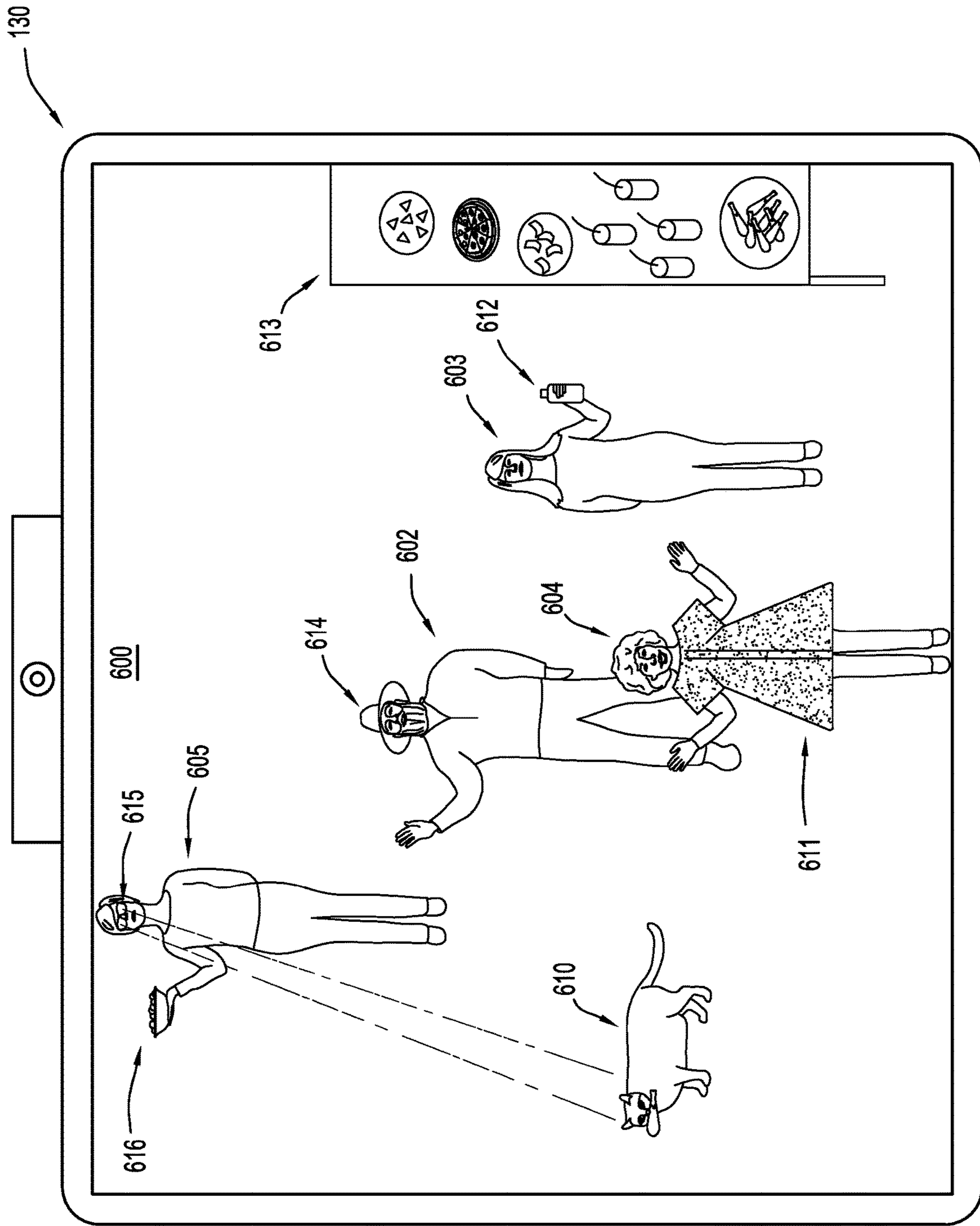


FIG. 6A

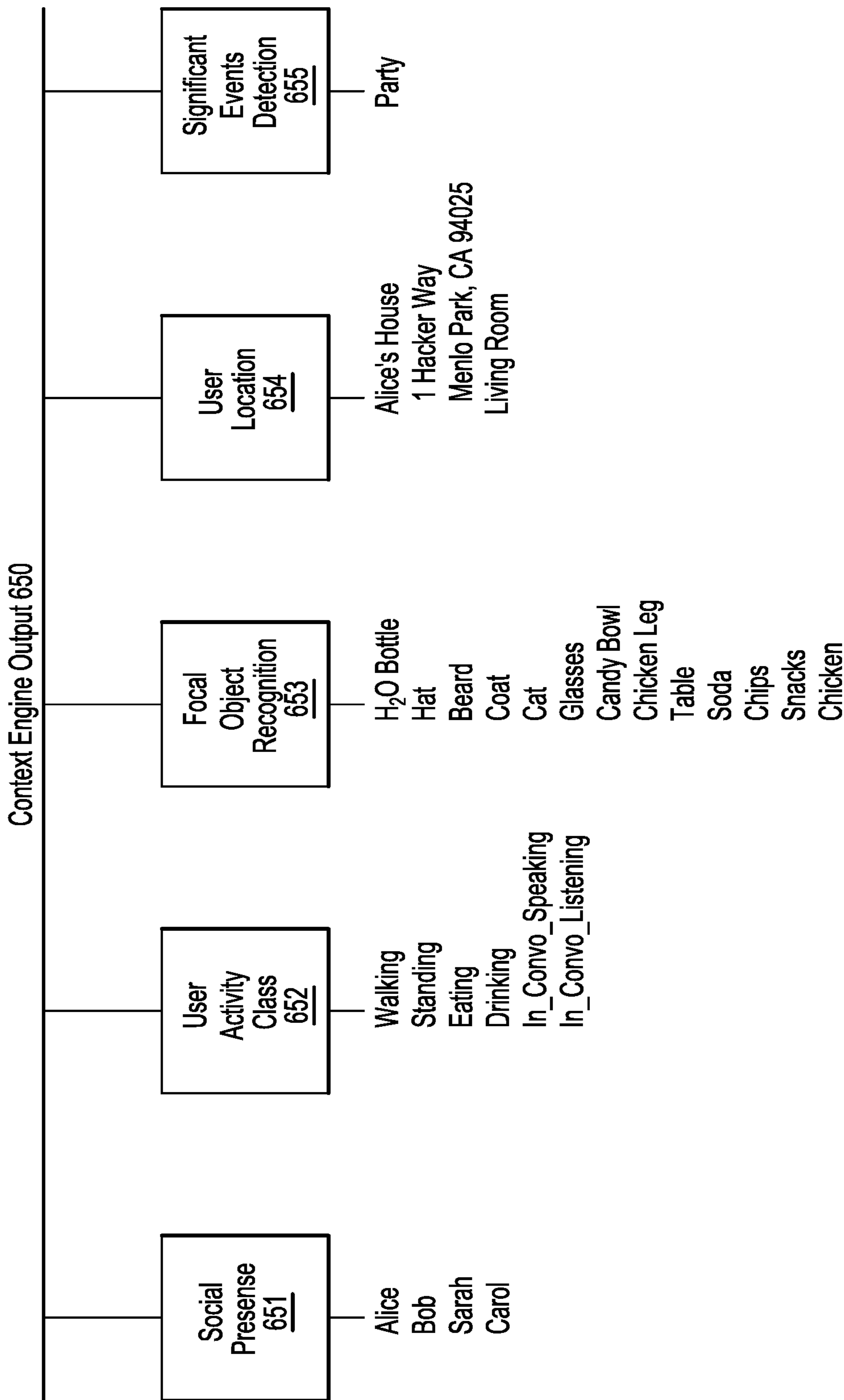


FIG. 6B

Scene Understanding Engine Output 660

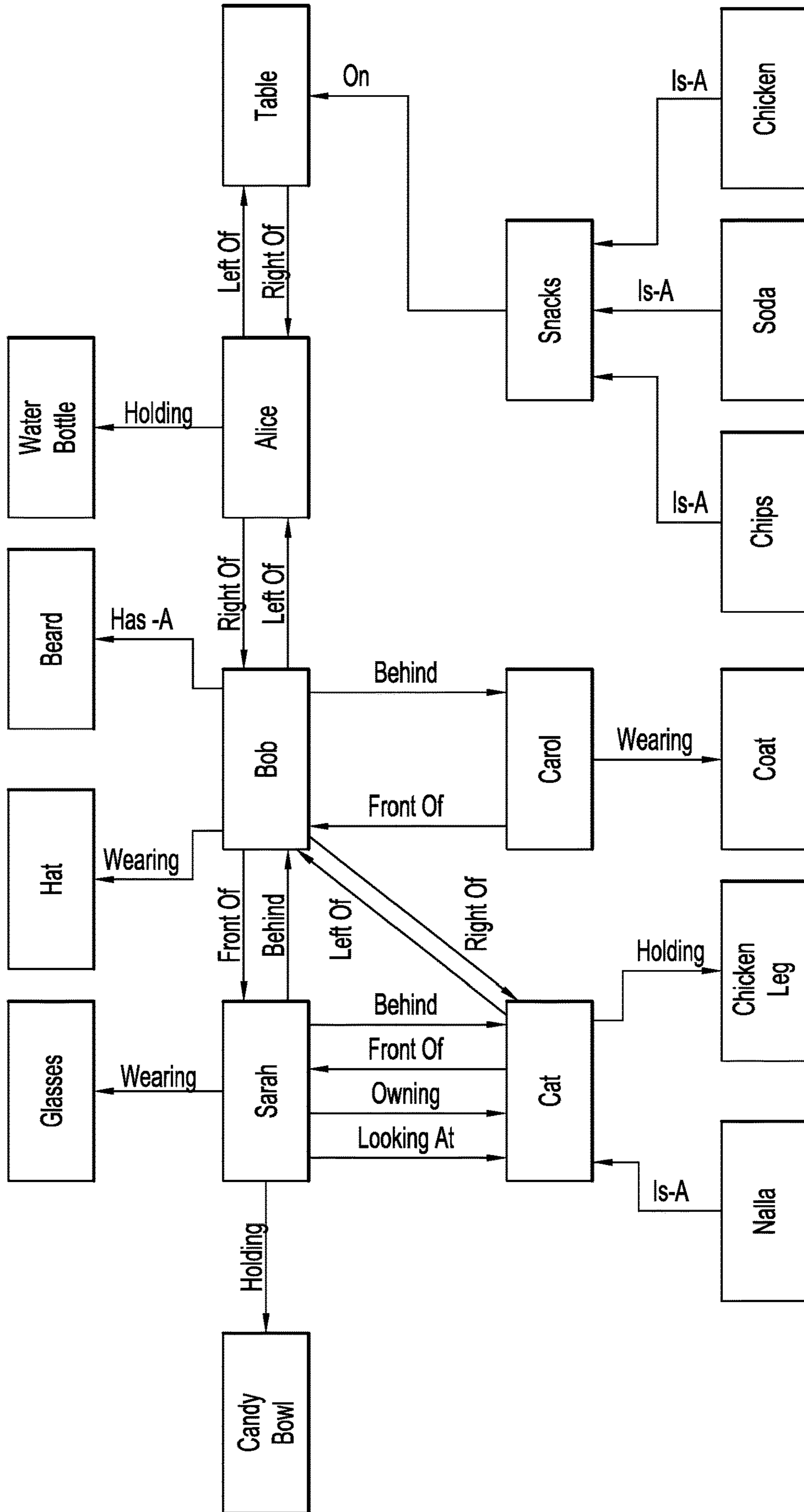


FIG. 6C

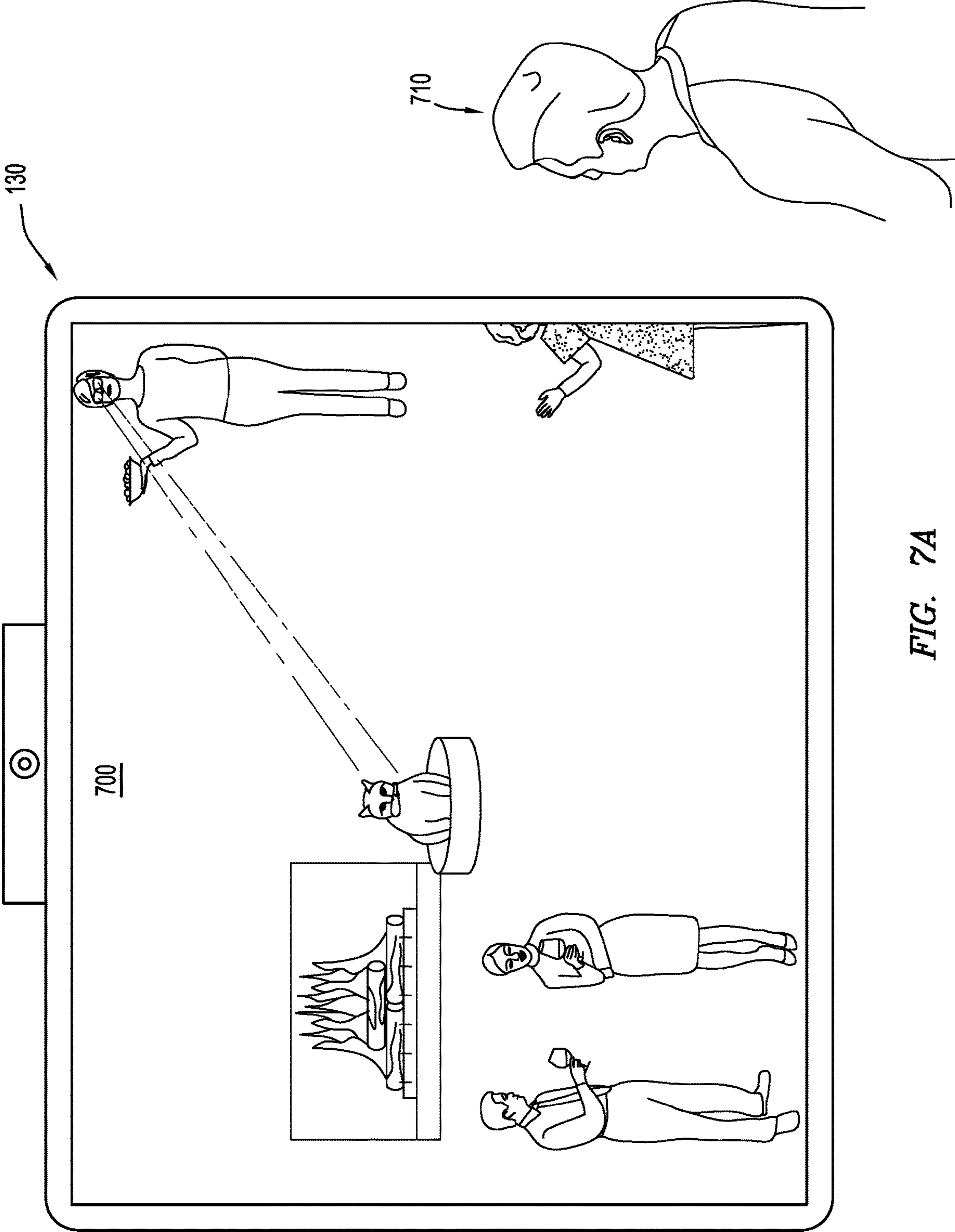


FIG. 7A

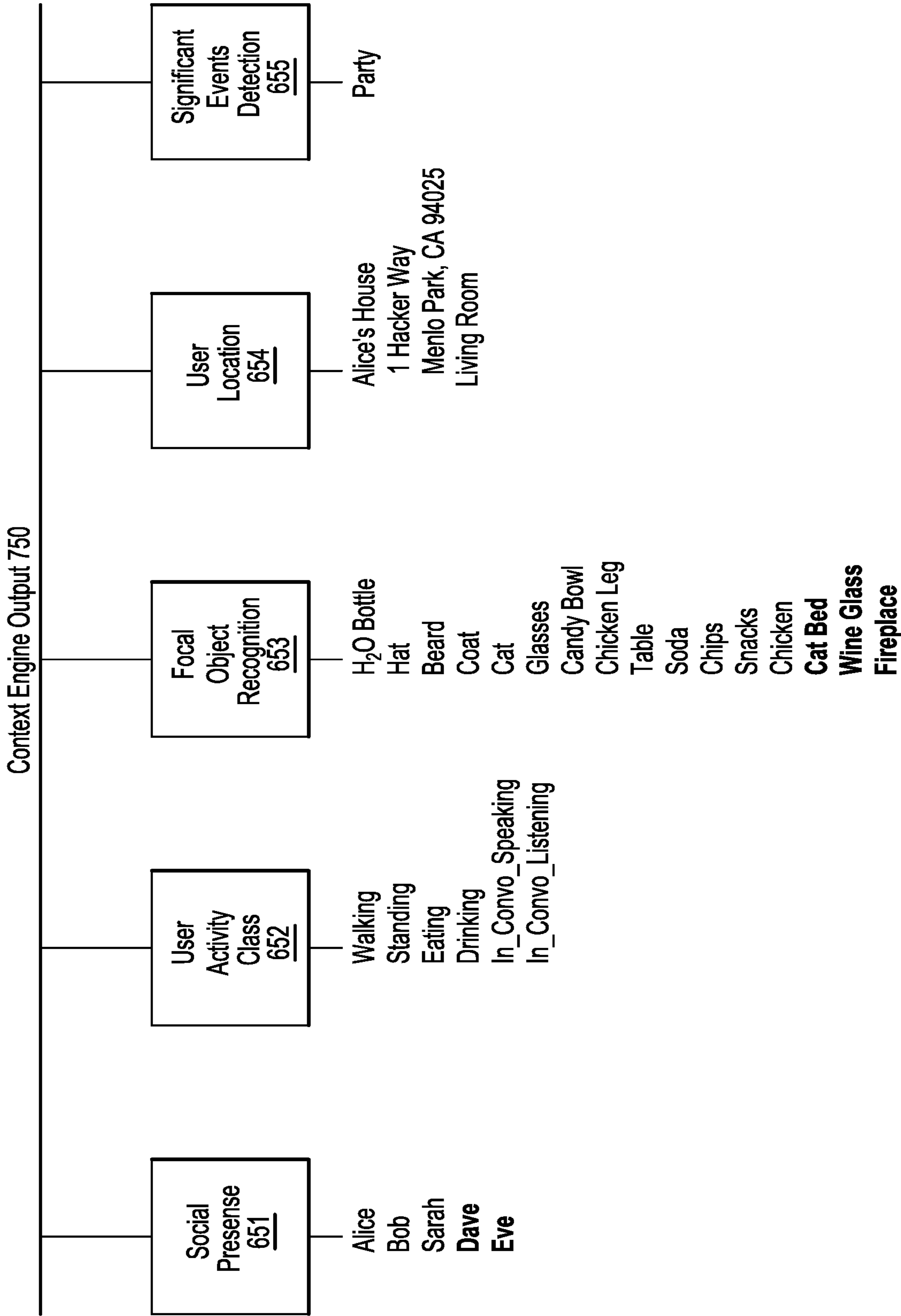


FIG. 7B

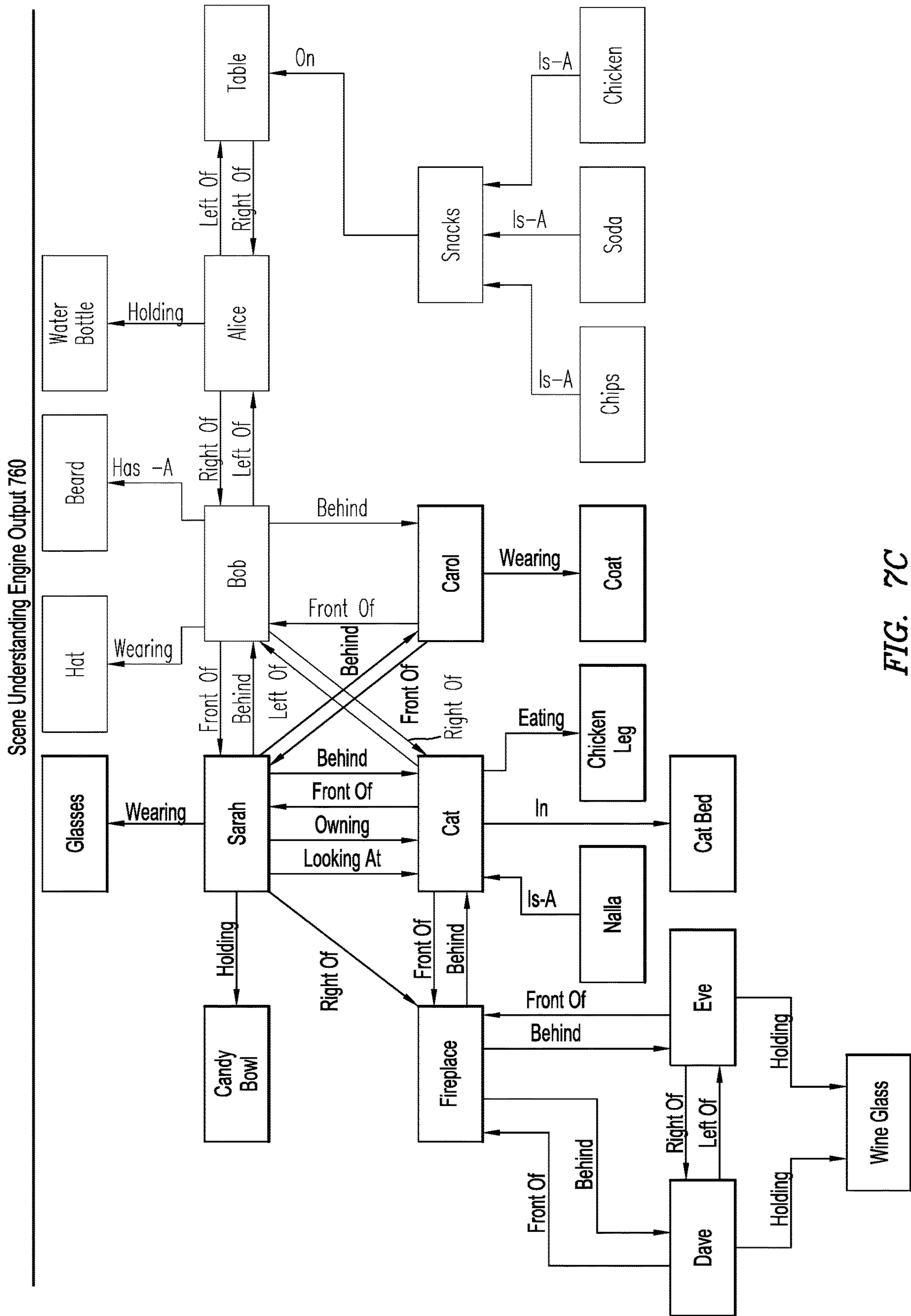


FIG. 7C

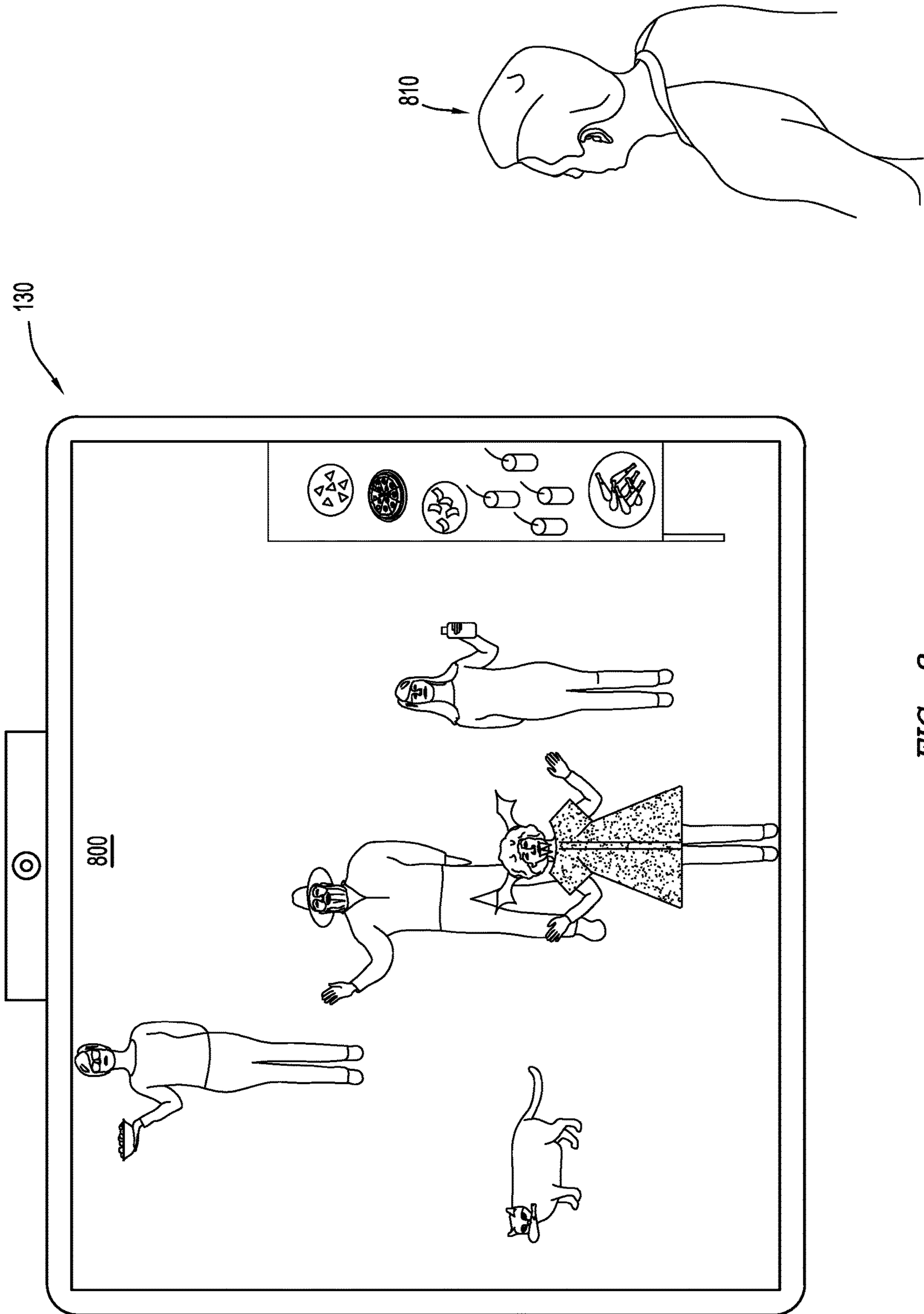


FIG. 8

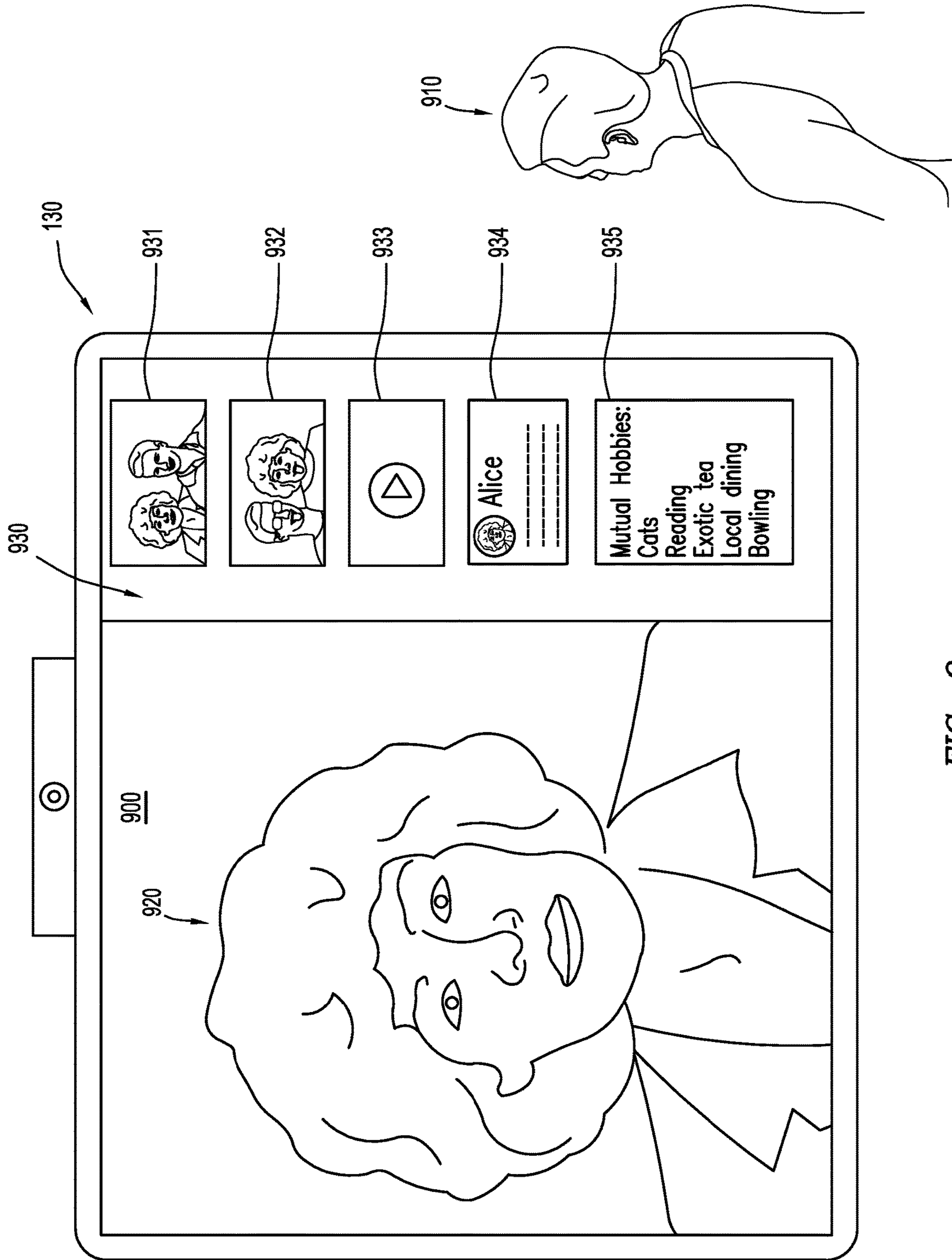
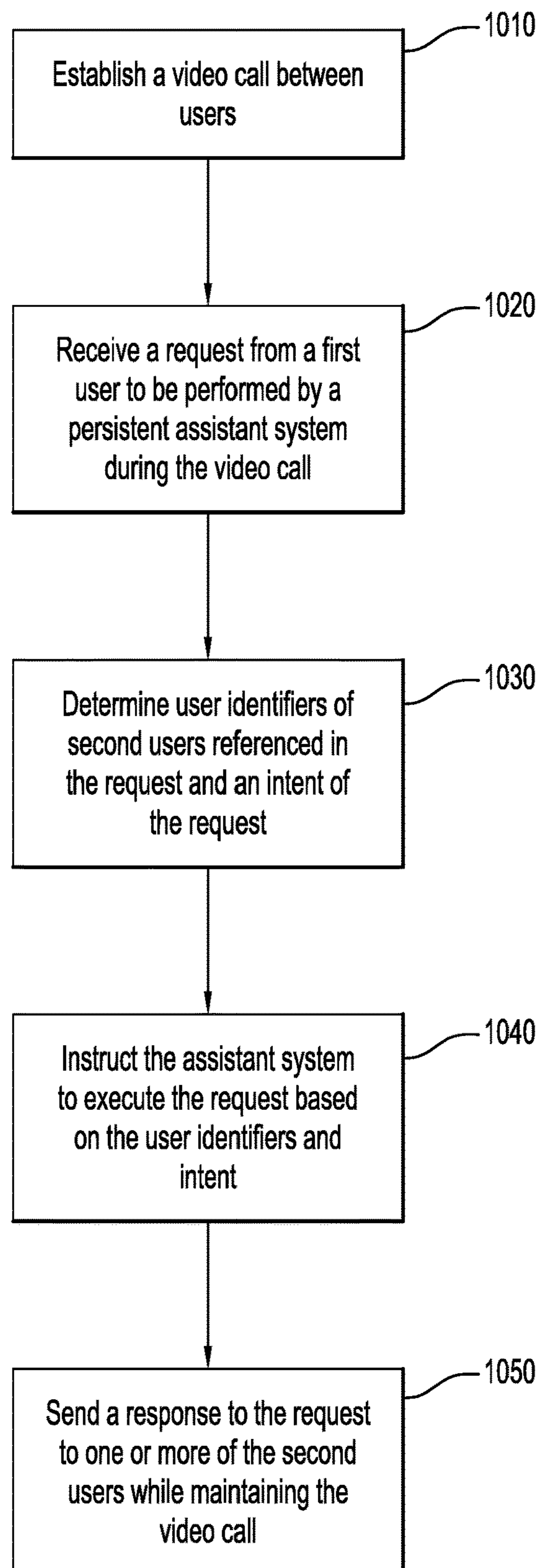


FIG. 9

1000*FIG. 10*

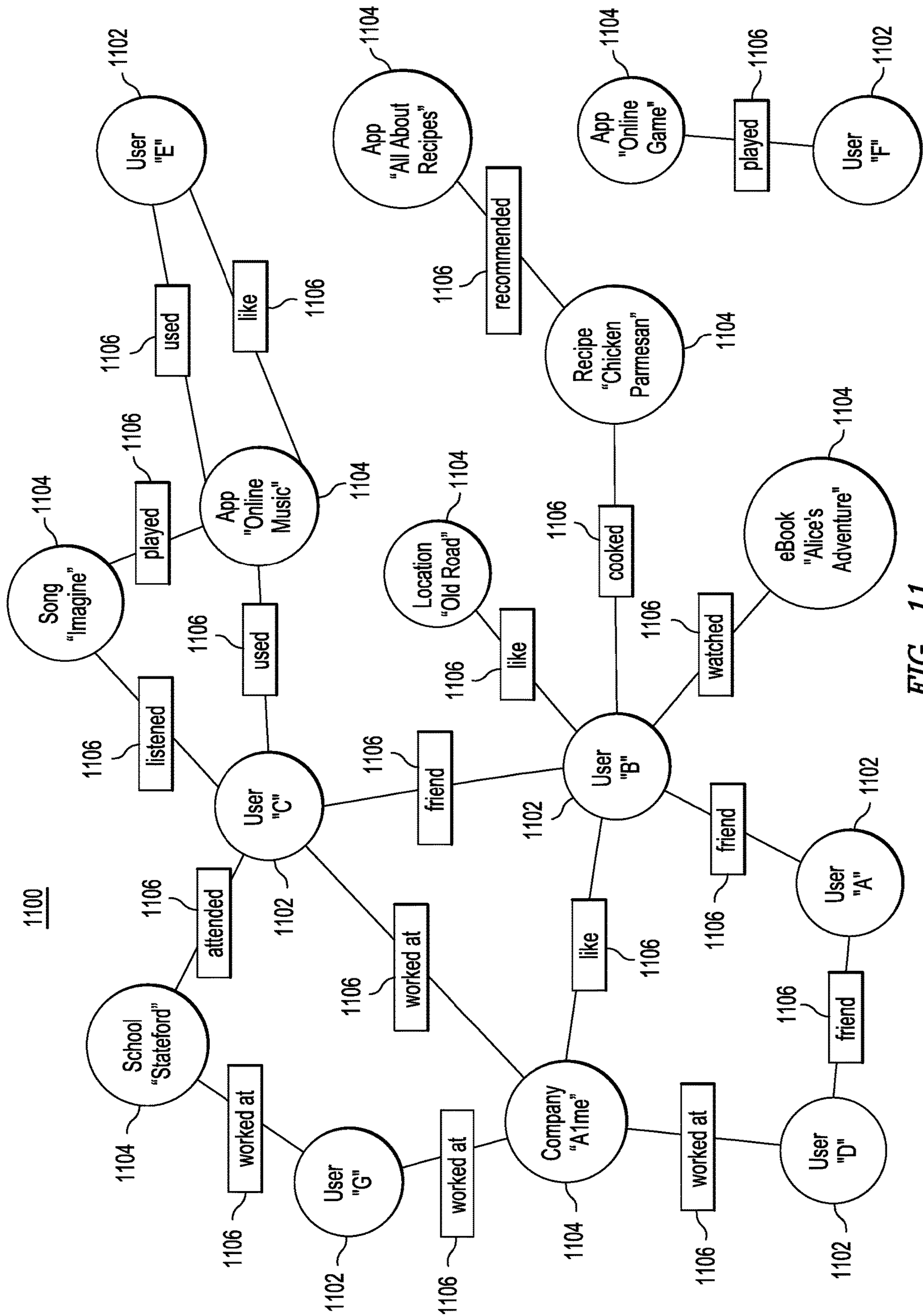


FIG. 11

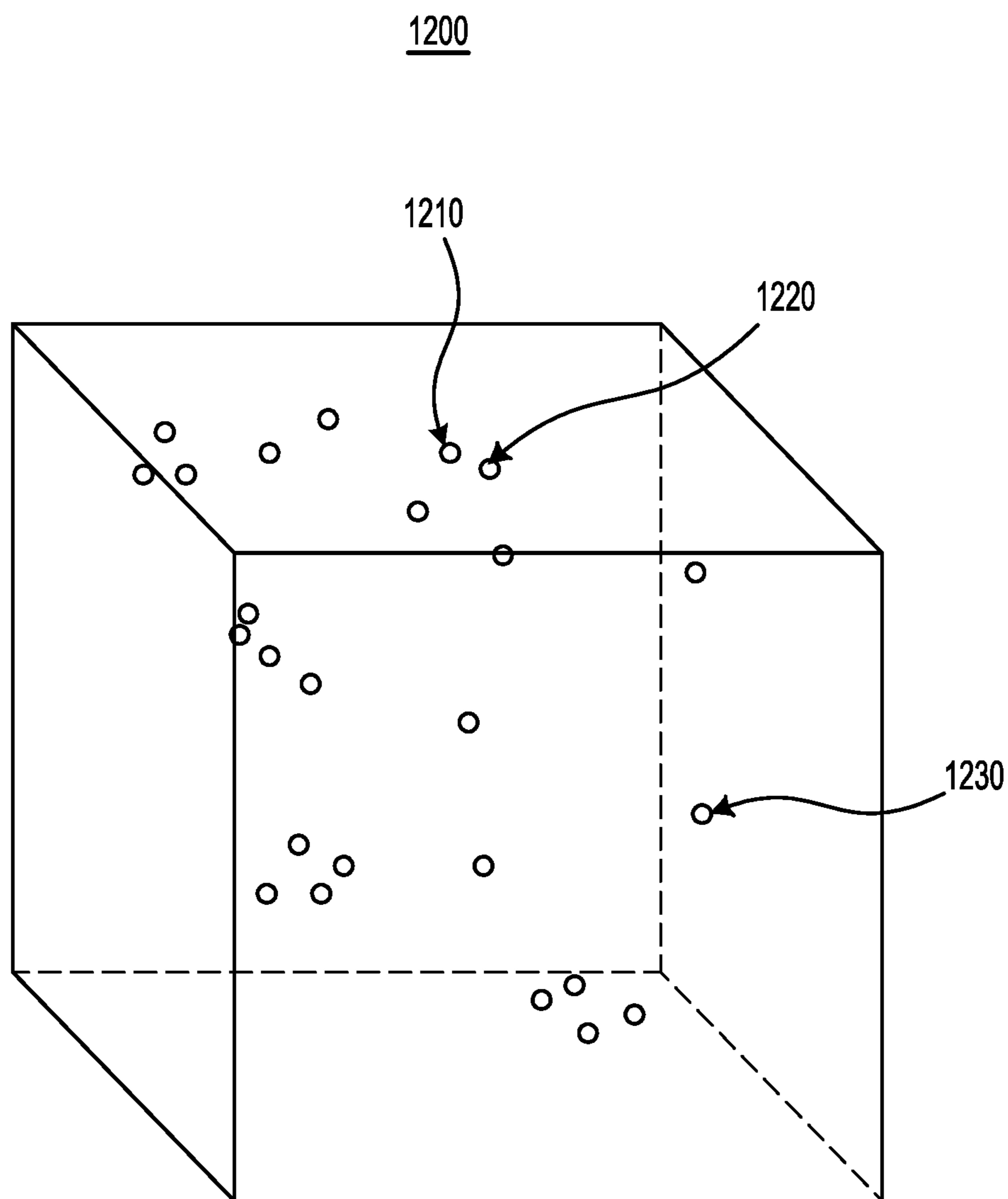


FIG. 12

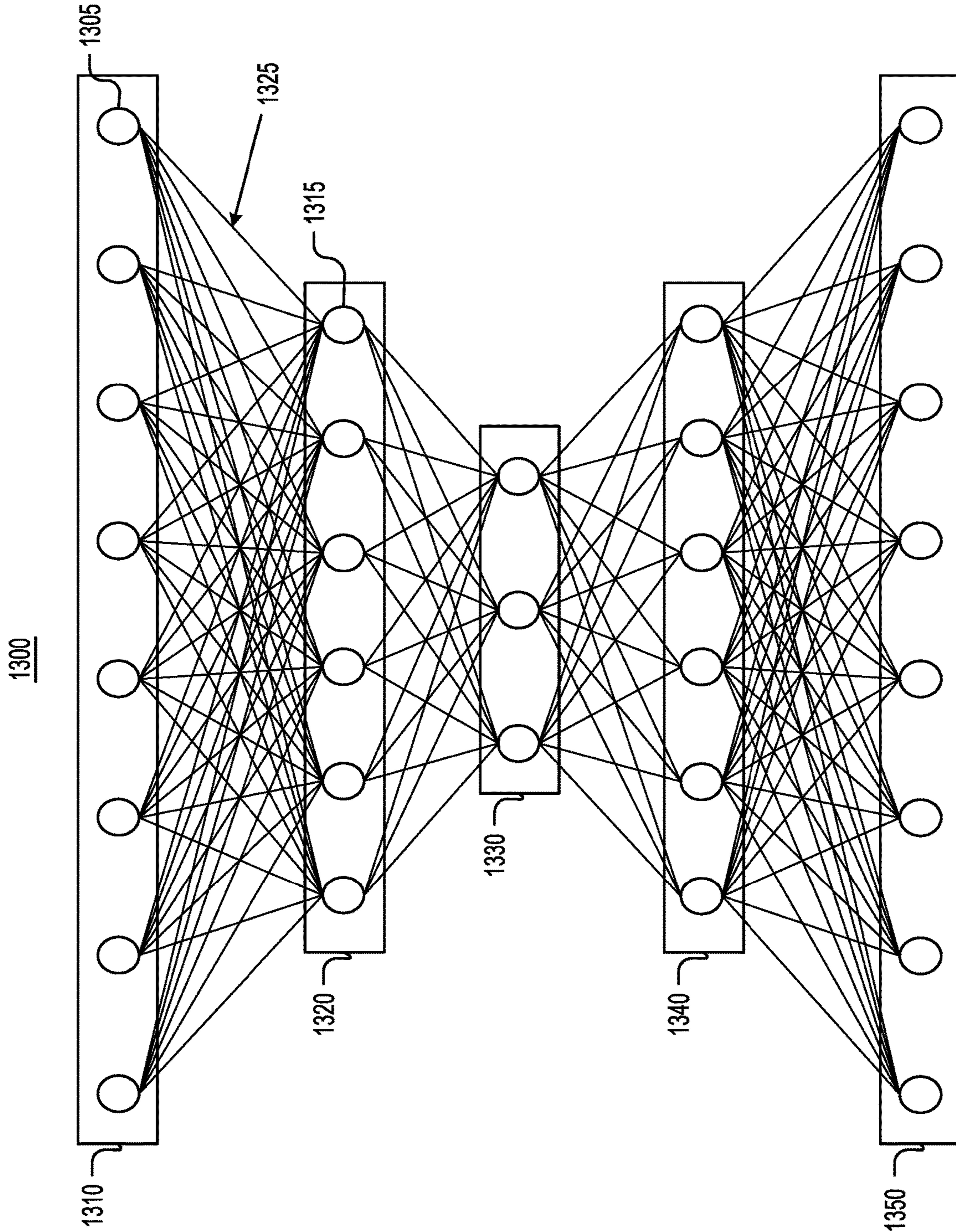


FIG. 13

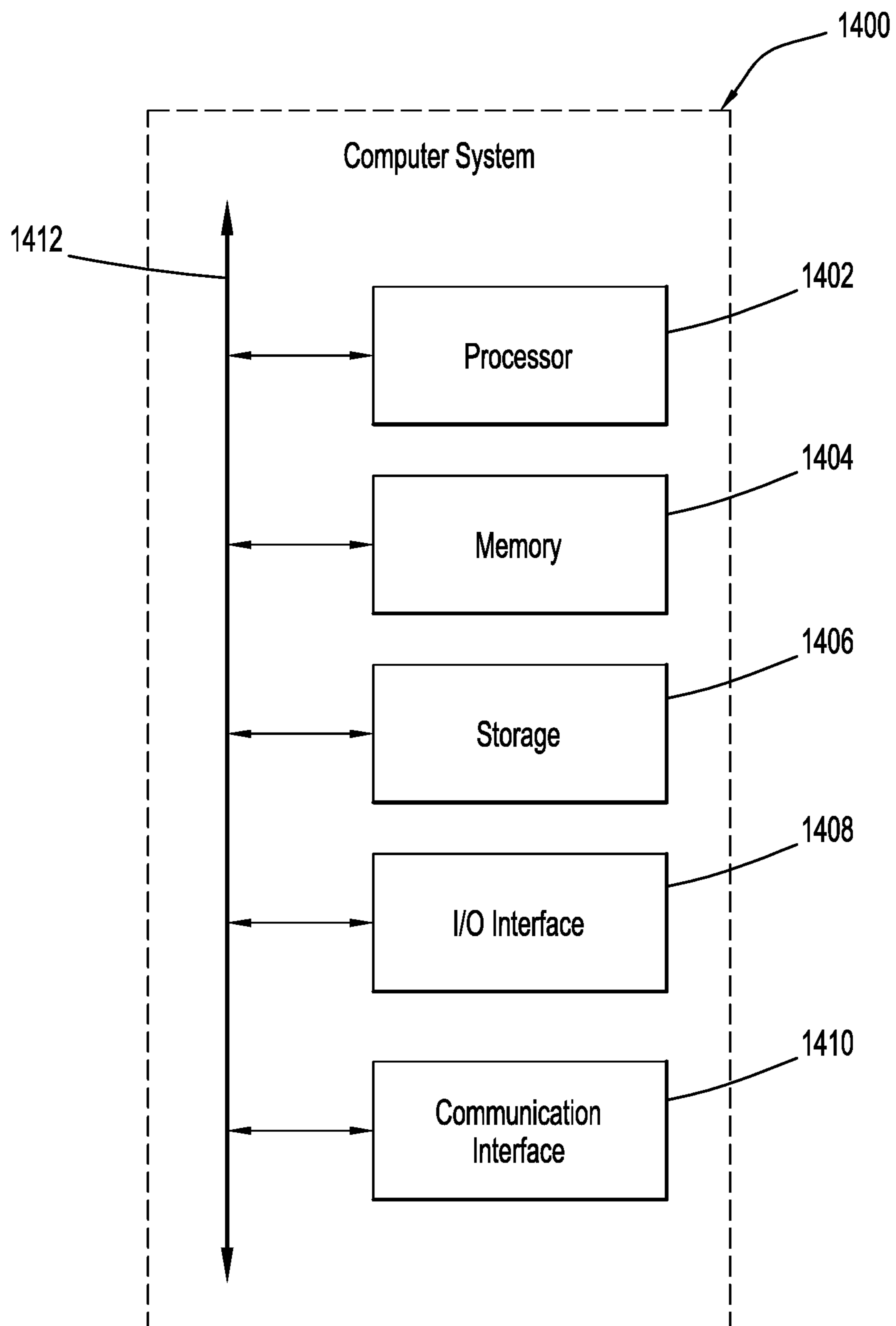


FIG. 14

IN-CALL EXPERIENCE ENHANCEMENT FOR ASSISTANT SYSTEMS

PRIORITY

This application claims the benefit under 35 U.S.C. § 119(e) of U.S. Provisional Patent Application No. 62/923,342, filed 18 Oct. 2019, which is incorporated herein by reference.

TECHNICAL FIELD

This disclosure generally relates to databases and file management within network environments, and in particular relates to hardware and software for smart assistant systems.

BACKGROUND

An assistant system can provide information or services on behalf of a user based on a combination of user input, location awareness, and the ability to access information from a variety of online sources (such as weather conditions, traffic congestion, news, stock prices, user schedules, retail prices, etc.). The user input may include text (e.g., online chat), especially in an instant messaging application or other applications, voice, images, motion, or a combination of them. The assistant system may perform concierge-type services (e.g., making dinner reservations, purchasing event tickets, making travel arrangements) or provide information based on the user input. The assistant system may also perform management or data-handling tasks based on online information and events without user initiation or interaction. Examples of those tasks that may be performed by an assistant system may include schedule management (e.g., sending an alert to a dinner date that a user is running late due to traffic conditions, update schedules for both parties, and change the restaurant reservation time). The assistant system may be enabled by the combination of computing devices, application programming interfaces (APIs), and the proliferation of applications on user devices.

A social-networking system, which may include a social-networking website, may enable its users (such as persons or organizations) to interact with it and with each other through it. The social-networking system may, with input from a user, create and store in the social-networking system a user profile associated with the user. The user profile may include demographic information, communication-channel information, and information on personal interests of the user. The social-networking system may also, with input from a user, create and store a record of relationships of the user with other users of the social-networking system, as well as provide services (e.g. profile/news feed posts, photo-sharing, event organization, messaging, games, or advertisements) to facilitate social interaction between or among users.

The social-networking system may send over one or more networks content or messages related to its services to a mobile or other computing device of a user. A user may also install software applications on a mobile or other computing device of the user for accessing a user profile of the user and other data within the social-networking system. The social-networking system may generate a personalized set of content objects to display to a user, such as a newsfeed of aggregated stories of other users connected to the user.

SUMMARY OF PARTICULAR EMBODIMENTS

In particular embodiments, the assistant system may assist a user to obtain information or services. The assistant system

may enable the user to interact with it with multi-modal user input (such as voice, text, image, video, motion) in stateful and multi-turn conversations to get assistance. As an example and not by way of limitation, the assistant system may support both audio (verbal) input and nonverbal input, such as vision, location, gesture, motion, or hybrid/multi-modal input. The assistant system may create and store a user profile comprising both personal and contextual information associated with the user. In particular embodiments, the assistant system may analyze the user input using natural-language understanding. The analysis may be based on the user profile of the user for more personalized and context-aware understanding. The assistant system may resolve entities associated with the user input based on the analysis. In particular embodiments, the assistant system may interact with different agents to obtain information or services that are associated with the resolved entities. The assistant system may generate a response for the user regarding the information or services by using natural-language generation. Through the interaction with the user, the assistant system may use dialog-management techniques to manage and advance the conversation flow with the user. In particular embodiments, the assistant system may further assist the user to effectively and efficiently digest the obtained information by summarizing the information. The assistant system may also assist the user to be more engaging with an online social network by providing tools that help the user interact with the online social network (e.g., creating posts, comments, messages). The assistant system may additionally assist the user to manage different tasks such as keeping track of events. In particular embodiments, the assistant system may proactively execute, without a user input, tasks that are relevant to user interests and preferences based on the user profile, at a time relevant for the user. In particular embodiments, the assistant system may check privacy settings to ensure that accessing a user's profile or other user information and executing different tasks are permitted subject to the user's privacy settings.

In particular embodiments, the assistant system may assist the user via a hybrid architecture built upon both client-side processes and server-side processes. The client-side processes and the server-side processes may be two parallel workflows for processing a user input and providing assistance to the user. In particular embodiments, the client-side processes may be performed locally on a client system associated with a user. By contrast, the server-side processes may be performed remotely on one or more computing systems. In particular embodiments, an arbitrator on the client system may coordinate receiving user input (e.g., an audio signal), determine whether to use a client-side process, a server-side process, or both, to respond to the user input, and analyze the processing results from each process. The arbitrator may instruct agents on the client-side or server-side to execute tasks associated with the user input based on the aforementioned analyses. The execution results may be further rendered as output to the client system. By leveraging both client-side and server-side processes, the assistant system can effectively assist a user with optimal usage of computing resources while at the same time protecting user privacy and enhancing security.

In particular embodiment, an in-call experience enhancement in which the assistant system is persistently active, but on standby during a call (such as a video or audio call) or other communication session (such as a text message thread), is provided. Such a persistently active assistant system may enable a user to invoke it in real-time during the call to execute tasks related to one or more other users on the

call. Furthermore, the persistently active assistant system may allow a single communication domain to be used in which the user can communicate with both other people via the call and with the assistant system itself. Current assistant systems typically go dormant during calls, so that a user must pause the call and reawaken the assistant system in order to issue commands. Thus, this single communication domain may greatly improve the user's experience, enabling a more social and natural interaction. The persistent assistant system may utilize an underlying multimodal architecture having separate context and scene understanding engines. The context engine may also be persistent during the call, gathering data for use by other modules in the assistant system that responds to a user query (subject to privacy settings). By contrast, the scene understanding engine may be awakened as needed to receive the data gathered by the context engine and determines a relationship among detected entities. Accordingly, with a video call in particular serving as a social experience backdrop, this persistent assistant system may enable numerous social, utility, communication, and image processing functionalities to be performed.

In particular embodiments, a video call between a plurality of client systems may be established, while persistently maintaining access to an assistant system during the video call. A request to be performed by the assistant system during the video call may then be received from a first client system of a first user. This request may reference one or more second users associated with second client systems. An intent of the request and one or more user identifiers of these one or more second users referenced by the request may be determined, and the assistant system may be instructed to execute the request based on the determined intent and user identifiers. Finally, a response to the request may be sent to one more of the plurality of client systems while maintaining the video call between the plurality of client systems.

Certain technical challenges exist in maintaining a quality video call between users. Video calls may lack a feeling of genuine social interaction; providing more social functions that may be performed during an actual video call may thus increase user interaction and satisfaction with the video call. However, one technical challenge to this may include identifying users in the video call that a viewing user in the video call wants to perform some social function with, as well as actually understanding the scene and context of the video call in order to more accurately execute the social function. A solution presented by embodiments disclosed herein to address this challenge may thus include continuously gathering context of the video call via a context engine and feeding this gathered information into a scene understanding engine, in order to generate relationship information between people and objects in the scene of the video call. Another technical challenge may be that, when conducting a video call on a client device, the user of that device may wish to preserve access to the functions of the device and access to a smart assistant system, which may go dormant during the video call. A solution presented by embodiments disclosed herein to address this challenge may thus involve a persistent assistant system that, rather than going dormant during a video call, remains active but on standby, and is thus accessible to the user to be invoked during a video call to execute various commands.

Certain embodiments disclosed herein may provide one or more technical advantages. As an example, accurately identifying any users and objects in a video call, as well as their context and relationship information (subject to privacy settings), may enable a viewing user to perform a variety of

social functions with respect to entities in the video call, even when the viewing user communicates those functions ambiguously. As another example, providing a persistent, always-on assistant system may enable a user to continue to use their client device and smart assistant normally, even while conducting a video call. Certain embodiments disclosed herein may provide none, some, or all of the above technical advantages. One or more other technical advantages may be readily apparent to one skilled in the art in view of the figures, descriptions, and claims of the present disclosure.

The embodiments disclosed herein are only examples, and the scope of this disclosure is not limited to them. Particular embodiments may include all, some, or none of the components, elements, features, functions, operations, or steps of the embodiments disclosed herein. Embodiments according to the invention are in particular disclosed in the attached claims directed to a method, a storage medium, a system and a computer program product, wherein any feature mentioned in one claim category, e.g. method, can be claimed in another claim category, e.g. system, as well. The dependencies or references back in the attached claims are chosen for formal reasons only. However, any subject matter resulting from a deliberate reference back to any previous claims (in particular multiple dependencies) can be claimed as well, so that any combination of claims and the features thereof are disclosed and can be claimed regardless of the dependencies chosen in the attached claims. The subject-matter which can be claimed comprises not only the combinations of features as set out in the attached claims but also any other combination of features in the claims, wherein each feature mentioned in the claims can be combined with any other feature or combination of other features in the claims. Furthermore, any of the embodiments and features described or depicted herein can be claimed in a separate claim and/or in any combination with any embodiment or feature described or depicted herein or with any of the features of the attached claims.

BRIEF DESCRIPTION OF THE DRAWINGS

FIG. 1 illustrates an example network environment associated with an assistant system.

FIG. 2 illustrates an example architecture of the assistant system.

FIG. 3 illustrates an example diagram flow of server-side processes of the assistant system.

FIG. 4 illustrates an example diagram flow of processing a user input by the assistant system.

FIG. 5 illustrates an example multimodal architecture of the assistant system.

FIG. 6A illustrates an example initial scene viewed during a video call on a first client system of a first user.

FIG. 6B illustrates an example chart of information of the scene generated by an always-on context engine.

FIG. 6C illustrates an example knowledge graph of the scene generated by a scene understanding engine.

FIG. 7A illustrates an example shifted scene viewed after a user command concerning an entity of the initial scene on the first client system of the first user.

FIG. 7B illustrates an example updated chart of information of the shifted scene generated by the context engine.

FIG. 7C illustrates an example updated knowledge graph of the shifted scene generated by the scene understanding engine.

5

FIG. 8 illustrates an example updated scene viewed after a user command concerning an entity of a previous scene on the first client system of the first user.

FIG. 9 illustrates an example video call in which content relevant to the video call is viewed on the client system of a user.

FIG. 10 illustrates an example method for generating a response to a user request to a persistent assistant system made during a call.

FIG. 11 illustrates an example social graph.

FIG. 12 illustrates an example view of an embedding space.

FIG. 13 illustrates an example artificial neural network.

FIG. 14 illustrates an example computer system.

DESCRIPTION OF EXAMPLE EMBODIMENTS

System Overview

FIG. 1 illustrates an example network environment 100 associated with an assistant system. Network environment 100 includes a client system 130, an assistant system 140, a social-networking system 160, and a third-party system 170 connected to each other by a network 110. Although FIG. 1 illustrates a particular arrangement of a client system 130, an assistant system 140, a social-networking system 160, a third-party system 170, and a network 110, this disclosure contemplates any suitable arrangement of a client system 130, an assistant system 140, a social-networking system 160, a third-party system 170, and a network 110. As an example and not by way of limitation, two or more of a client system 130, a social-networking system 160, an assistant system 140, and a third-party system 170 may be connected to each other directly, bypassing a network 110. As another example, two or more of a client system 130, an assistant system 140, a social-networking system 160, and a third-party system 170 may be physically or logically co-located with each other in whole or in part. Moreover, although FIG. 1 illustrates a particular number of client systems 130, assistant systems 140, social-networking systems 160, third-party systems 170, and networks 110, this disclosure contemplates any suitable number of client systems 130, assistant systems 140, social-networking systems 160, third-party systems 170, and networks 110. As an example and not by way of limitation, network environment 100 may include multiple client systems 130, assistant systems 140, social-networking systems 160, third-party systems 170, and networks 110.

This disclosure contemplates any suitable network 110. As an example and not by way of limitation, one or more portions of a network 110 may include an ad hoc network, an intranet, an extranet, a virtual private network (VPN), a local area network (LAN), a wireless LAN (WLAN), a wide area network (WAN), a wireless WAN (WWAN), a metropolitan area network (MAN), a portion of the Internet, a portion of the Public Switched Telephone Network (PSTN), a cellular telephone network, or a combination of two or more of these. A network 110 may include one or more networks 110.

Links 150 may connect a client system 130, an assistant system 140, a social-networking system 160, and a third-party system 170 to a communication network 110 or to each other. This disclosure contemplates any suitable links 150. In particular embodiments, one or more links 150 include one or more wireline (such as for example Digital Subscriber Line (DSL) or Data Over Cable Service Interface Specification (DOCSIS)), wireless (such as for example Wi-Fi or Worldwide Interoperability for Microwave Access

6

(WiMAX)), or optical (such as for example Synchronous Optical Network (SONET) or Synchronous Digital Hierarchy (SDH)) links. In particular embodiments, one or more links 150 each include an ad hoc network, an intranet, an extranet, a VPN, a LAN, a WLAN, a WAN, a WWAN, a MAN, a portion of the Internet, a portion of the PSTN, a cellular technology-based network, a satellite communications technology-based network, another link 150, or a combination of two or more such links 150. Links 150 need not necessarily be the same throughout a network environment 100. One or more first links 150 may differ in one or more respects from one or more second links 150.

In particular embodiments, a client system 130 may be an electronic device including hardware, software, or embedded logic components or a combination of two or more such components and capable of carrying out the appropriate functionalities implemented or supported by a client system 130. As an example and not by way of limitation, a client system 130 may include a computer system such as a desktop computer, notebook or laptop computer, netbook, a tablet computer, e-book reader, GPS device, camera, personal digital assistant (PDA), handheld electronic device, cellular telephone, smartphone, smart speaker, virtual reality (VR) headset, augment reality (AR) smart glasses, other suitable electronic device, or any suitable combination thereof. In particular embodiments, the client system 130 may be a smart assistant device. More information on smart assistant devices may be found in U.S. patent application Ser. No. 15/949,011, filed 9 Apr. 2018, U.S. patent application Ser. No. 16/153,574, filed 5 Oct. 2018, U.S. Design patent application Ser. No. 29/631,910, filed 3 Jan. 2018, U.S. Design patent application Ser. No. 29/631,747, filed 2 Jan. 2018, U.S. Design patent application Ser. No. 29/631,913, filed 3 Jan. 2018, and U.S. Design patent application Ser. No. 29/631,914, filed 3 Jan. 2018, each of which is incorporated by reference. This disclosure contemplates any suitable client systems 130. A client system 130 may enable a network user at a client system 130 to access a network 110. A client system 130 may enable its user to communicate with other users at other client systems 130.

In particular embodiments, a client system 130 may include a web browser 132, and may have one or more add-ons, plug-ins, or other extensions. A user at a client system 130 may enter a Uniform Resource Locator (URL) or other address directing a web browser 132 to a particular server (such as server 162, or a server associated with a third-party system 170), and the web browser 132 may generate a Hyper Text Transfer Protocol (HTTP) request and communicate the HTTP request to server. The server may accept the HTTP request and communicate to a client system 130 one or more Hyper Text Markup Language (HTML) files responsive to the HTTP request. The client system 130 may render a web interface (e.g. a webpage) based on the HTML files from the server for presentation to the user. This disclosure contemplates any suitable source files. As an example and not by way of limitation, a web interface may be rendered from HTML files, Extensible Hyper Text Markup Language (XHTML) files, or Extensible Markup Language (XML) files, according to particular needs. Such interfaces may also execute scripts, combinations of markup language and scripts, and the like. Herein, reference to a web interface encompasses one or more corresponding source files (which a browser may use to render the web interface) and vice versa, where appropriate.

In particular embodiments, a client system 130 may include a social-networking application 134 installed on the client system 130. A user at a client system 130 may use the

social-networking application **134** to access on online social network. The user at the client system **130** may use the social-networking application **134** to communicate with the user's social connections (e.g., friends, followers, followed accounts, contacts, etc.). The user at the client system **130** may also use the social-networking application **134** to interact with a plurality of content objects (e.g., posts, news articles, ephemeral content, etc.) on the online social network. As an example and not by way of limitation, the user may browse trending topics and breaking news using the social-networking application **134**.

In particular embodiments, a client system **130** may include an assistant application **136**. A user at a client system **130** may use the assistant application **136** to interact with the assistant system **140**. In particular embodiments, the assistant application **136** may comprise a stand-alone application. In particular embodiments, the assistant application **136** may be integrated into the social-networking application **134** or another suitable application (e.g., a messaging application). In particular embodiments, the assistant application **136** may be also integrated into the client system **130**, an assistant hardware device, or any other suitable hardware devices. In particular embodiments, the assistant application **136** may be accessed via the web browser **132**. In particular embodiments, the user may provide input via different modalities. As an example and not by way of limitation, the modalities may include audio, text, image, video, motion, orientation, etc. The assistant application **136** may communicate the user input to the assistant system **140**. Based on the user input, the assistant system **140** may generate responses. The assistant system **140** may send the generated responses to the assistant application **136**. The assistant application **136** may then present the responses to the user at the client system **130**. The presented responses may be based on different modalities such as audio, text, image, and video. As an example and not by way of limitation, the user may verbally ask the assistant application **136** about the traffic information (i.e., via an audio modality) by speaking into a microphone of the client system **130**. The assistant application **136** may then communicate the request to the assistant system **140**. The assistant system **140** may accordingly generate a response and send it back to the assistant application **136**. The assistant application **136** may further present the response to the user in text and/or images on a display of the client system **130**.

In particular embodiments, an assistant system **140** may assist users to retrieve information from different sources. The assistant system **140** may also assist user to request services from different service providers. In particular embodiments, the assist system **140** may receive a user request for information or services via the assistant application **136** in the client system **130**. The assist system **140** may use natural-language understanding to analyze the user request based on user's profile and other relevant information. The result of the analysis may comprise different entities associated with an online social network. The assistant system **140** may then retrieve information or request services associated with these entities. In particular embodiments, the assistant system **140** may interact with the social-networking system **160** and/or third-party system **170** when retrieving information or requesting services for the user. In particular embodiments, the assistant system **140** may generate a personalized communication content for the user using natural-language generating techniques. The personalized communication content may comprise, for example, the retrieved information or the status of the requested services. In particular embodiments, the assistant

system **140** may enable the user to interact with it regarding the information or services in a stateful and multi-turn conversation by using dialog-management techniques. The functionality of the assistant system **140** is described in more detail in the discussion of FIG. 2 below.

In particular embodiments, the social-networking system **160** may be a network-addressable computing system that can host an online social network. The social-networking system **160** may generate, store, receive, and send social-networking data, such as, for example, user profile data, concept-profile data, social-graph information, or other suitable data related to the online social network. The social-networking system **160** may be accessed by the other components of network environment **100** either directly or via a network **110**. As an example and not by way of limitation, a client system **130** may access the social-networking system **160** using a web browser **132**, or a native application associated with the social-networking system **160** (e.g., a mobile social-networking application, a messaging application, another suitable application, or any combination thereof) either directly or via a network **110**. In particular embodiments, the social-networking system **160** may include one or more servers **162**. Each server **162** may be a unitary server or a distributed server spanning multiple computers or multiple datacenters. Servers **162** may be of various types, such as, for example and without limitation, web server, news server, mail server, message server, advertising server, file server, application server, exchange server, database server, proxy server, another server suitable for performing functions or processes described herein, or any combination thereof. In particular embodiments, each server **162** may include hardware, software, or embedded logic components or a combination of two or more such components for carrying out the appropriate functionalities implemented or supported by server **162**. In particular embodiments, the social-networking system **160** may include one or more data stores **164**. Data stores **164** may be used to store various types of information. In particular embodiments, the information stored in data stores **164** may be organized according to specific data structures. In particular embodiments, each data store **164** may be a relational, columnar, correlation, or other suitable database. Although this disclosure describes or illustrates particular types of databases, this disclosure contemplates any suitable types of databases. Particular embodiments may provide interfaces that enable a client system **130**, a social-networking system **160**, an assistant system **140**, or a third-party system **170** to manage, retrieve, modify, add, or delete, the information stored in data store **164**.

In particular embodiments, the social-networking system **160** may store one or more social graphs in one or more data stores **164**. In particular embodiments, a social graph may include multiple nodes—which may include multiple user nodes (each corresponding to a particular user) or multiple concept nodes (each corresponding to a particular concept)—and multiple edges connecting the nodes. The social-networking system **160** may provide users of the online social network the ability to communicate and interact with other users. In particular embodiments, users may join the online social network via the social-networking system **160** and then add connections (e.g., relationships) to a number of other users of the social-networking system **160** whom they want to be connected to. Herein, the term “friend” may refer to any other user of the social-networking system **160** with whom a user has formed a connection, association, or relationship via the social-networking system **160**.

In particular embodiments, the social-networking system **160** may provide users with the ability to take actions on various types of items or objects, supported by the social-networking system **160**. As an example and not by way of limitation, the items and objects may include groups or social networks to which users of the social-networking system **160** may belong, events or calendar entries in which a user might be interested, computer-based applications that a user may use, transactions that allow users to buy or sell items via the service, interactions with advertisements that a user may perform, or other suitable items or objects. A user may interact with anything that is capable of being represented in the social-networking system **160** or by an external system of a third-party system **170**, which is separate from the social-networking system **160** and coupled to the social-networking system **160** via a network **110**.

In particular embodiments, the social-networking system **160** may be capable of linking a variety of entities. As an example and not by way of limitation, the social-networking system **160** may enable users to interact with each other as well as receive content from third-party systems **170** or other entities, or to allow users to interact with these entities through an application programming interfaces (API) or other communication channels.

In particular embodiments, a third-party system **170** may include one or more types of servers, one or more data stores, one or more interfaces, including but not limited to APIs, one or more web services, one or more content sources, one or more networks, or any other suitable components, e.g., that servers may communicate with. A third-party system **170** may be operated by a different entity from an entity operating the social-networking system **160**. In particular embodiments, however, the social-networking system **160** and third-party systems **170** may operate in conjunction with each other to provide social-networking services to users of the social-networking system **160** or third-party systems **170**. In this sense, the social-networking system **160** may provide a platform, or backbone, which other systems, such as third-party systems **170**, may use to provide social-networking services and functionality to users across the Internet.

In particular embodiments, a third-party system **170** may include a third-party content object provider. A third-party content object provider may include one or more sources of content objects, which may be communicated to a client system **130**. As an example and not by way of limitation, content objects may include information regarding things or activities of interest to the user, such as, for example, movie show times, movie reviews, restaurant reviews, restaurant menus, product information and reviews, or other suitable information. As another example and not by way of limitation, content objects may include incentive content objects, such as coupons, discount tickets, gift certificates, or other suitable incentive objects. In particular embodiments, a third-party content provider may use one or more third-party agents to provide content objects and/or services. A third-party agent may be an implementation that is hosted and executing on the third-party system **170**.

In particular embodiments, the social-networking system **160** also includes user-generated content objects, which may enhance a user's interactions with the social-networking system **160**. User-generated content may include anything a user can add, upload, send, or "post" to the social-networking system **160**. As an example and not by way of limitation, a user communicates posts to the social-networking system **160** from a client system **130**. Posts may include data such as status updates or other textual data, location information,

photos, videos, links, music or other similar data or media. Content may also be added to the social-networking system **160** by a third-party through a "communication channel," such as a newsfeed or stream.

In particular embodiments, the social-networking system **160** may include a variety of servers, sub-systems, programs, modules, logs, and data stores. In particular embodiments, the social-networking system **160** may include one or more of the following: a web server, action logger, API-request server, relevance-and-ranking engine, content-object classifier, notification controller, action log, third-party-content-object-exposure log, inference module, authorization/privacy server, search module, advertisement-targeting module, user-interface module, user-profile store, connection store, third-party content store, or location store. The social-networking system **160** may also include suitable components such as network interfaces, security mechanisms, load balancers, failover servers, management-and-network-operations consoles, other suitable components, or any suitable combination thereof. In particular embodiments, the social-networking system **160** may include one or more user-profile stores for storing user profiles. A user profile may include, for example, biographic information, demographic information, behavioral information, social information, or other types of descriptive information, such as work experience, educational history, hobbies or preferences, interests, affinities, or location. Interest information may include interests related to one or more categories. Categories may be general or specific. As an example and not by way of limitation, if a user "likes" an article about a brand of shoes the category may be the brand, or the general category of "shoes" or "clothing." A connection store may be used for storing connection information about users. The connection information may indicate users who have similar or common work experience, group memberships, hobbies, educational history, or are in any way related or share common attributes. The connection information may also include user-defined connections between different users and content (both internal and external). A web server may be used for linking the social-networking system **160** to one or more client systems **130** or one or more third-party systems **170** via a network **110**. The web server may include a mail server or other messaging functionality for receiving and routing messages between the social-networking system **160** and one or more client systems **130**. An API-request server may allow, for example, an assistant system **140** or a third-party system **170** to access information from the social-networking system **160** by calling one or more APIs. An action logger may be used to receive communications from a web server about a user's actions on or off the social-networking system **160**. In conjunction with the action log, a third-party-content-object log may be maintained of user exposures to third-party-content objects. A notification controller may provide information regarding content objects to a client system **130**. Information may be pushed to a client system **130** as notifications, or information may be pulled from a client system **130** responsive to a request received from a client system **130**. Authorization servers may be used to enforce one or more privacy settings of the users of the social-networking system **160**. A privacy setting of a user determines how particular information associated with a user can be shared. The authorization server may allow users to opt in to or opt out of having their actions logged by the social-networking system **160** or shared with other systems (e.g., a third-party system **170**), such as, for example, by setting appropriate privacy settings. Third-party-content-object stores may be used to store con-

tent objects received from third parties, such as a third-party system 170. Location stores may be used for storing location information received from client systems 130 associated with users. Advertisement-pricing modules may combine social information, the current time, location information, or other suitable information to provide relevant advertisements, in the form of notifications, to a user.

Assistant Systems

FIG. 2 illustrates an example architecture of an assistant system 140. In particular embodiments, the assistant system 140 may assist a user to obtain information or services. The assistant system 140 may enable the user to interact with it with multi-modal user input (such as voice, text, image, video, motion) in stateful and multi-turn conversations to get assistance. As an example and not by way of limitation, the assistant system 140 may support both audio input (verbal) and nonverbal input, such as vision, location, gesture, motion, or hybrid/multi-modal input. The assistant system 140 may create and store a user profile comprising both personal and contextual information associated with the user. In particular embodiments, the assistant system 140 may analyze the user input using natural-language understanding. The analysis may be based on the user profile of the user for more personalized and context-aware understanding. The assistant system 140 may resolve entities associated with the user input based on the analysis. In particular embodiments, the assistant system 140 may interact with different agents to obtain information or services that are associated with the resolved entities. The assistant system 140 may generate a response for the user regarding the information or services by using natural-language generation. Through the interaction with the user, the assistant system 140 may use dialog management techniques to manage and forward the conversation flow with the user. In particular embodiments, the assistant system 140 may further assist the user to effectively and efficiently digest the obtained information by summarizing the information. The assistant system 140 may also assist the user to be more engaging with an online social network by providing tools that help the user interact with the online social network (e.g., creating posts, comments, messages). The assistant system 140 may additionally assist the user to manage different tasks such as keeping track of events. In particular embodiments, the assistant system 140 may proactively execute, without a user input, pre-authorized tasks that are relevant to user interests and preferences based on the user profile, at a time relevant for the user. In particular embodiments, the assistant system 140 may check privacy settings to ensure that accessing a user's profile or other user information and executing different tasks are permitted subject to the user's privacy settings. More information on assisting users subject to privacy settings may be found in U.S. patent application Ser. No. 16/182,542, filed 6 Nov. 2018, which is incorporated by reference.

In particular embodiments, the assistant system 140 may assist the user via a hybrid architecture built upon both client-side processes and server-side processes. The client-side processes and the server-side processes may be two parallel workflows for processing a user input and providing assistances to the user. In particular embodiments, the client-side processes may be performed locally on a client system 130 associated with a user. By contrast, the server-side processes may be performed remotely on one or more computing systems. In particular embodiments, an assistant orchestrator on the client system 130 may coordinate receiving user input (e.g., audio signal) and determine whether to use client-side processes, server-side processes, or both, to

respond to the user input. A dialog arbitrator may analyze the processing results from each process. The dialog arbitrator may instruct agents on the client-side or server-side to execute tasks associated with the user input based on the aforementioned analyses. The execution results may be further rendered as output to the client system 130. By leveraging both client-side and server-side processes, the assistant system 140 can effectively assist a user with optimal usage of computing resources while at the same time protecting user privacy and enhancing security.

In particular embodiments, the assistant system 140 may receive a user input from a client system 130 associated with the user. In particular embodiments, the user input may be a user-generated input that is sent to the assistant system 140 in a single turn. The user input may be verbal, nonverbal, or a combination thereof. As an example and not by way of limitation, the nonverbal user input may be based on the user's voice, vision, location, activity, gesture, motion, or a combination thereof. If the user input is based on the user's voice (e.g., the user may speak to the client system 130), such user input may be first processed by a system audio API 202 (application programming interface). The system audio API 202 may conduct echo cancellation, noise removal, beam forming, and self-user voice activation, speaker identification, voice activity detection (VAD), and any other acoustic techniques to generate audio data that is readily processable by the assistant system 140. In particular embodiments, the system audio API 202 may perform wake-word detection 204 from the user input. As an example and not by way of limitation, a wake-word may be "hey assistant". If such wake-word is detected, the assistant system 140 may be activated accordingly. In alternative embodiments, the user may activate the assistant system 140 via a visual signal without a wake-word. The visual signal may be received at a low-power sensor (e.g., a camera) that can detect various visual signals. As an example and not by way of limitation, the visual signal may be a barcode, a QR code or a universal product code (UPC) detected by the client system 130. As another example and not by way of limitation, the visual signal may be the user's gaze at an object. As yet another example and not by way of limitation, the visual signal may be a user gesture, e.g., the user pointing at an object.

In particular embodiments, the audio data from the system audio API 202 may be sent to an assistant orchestrator 206. The assistant orchestrator 206 may be executing on the client system 130. In particular embodiments, the assistant orchestrator 206 may determine whether to respond to the user input by using client-side processes, server-side processes, or both. As indicated in FIG. 2, the client-side processes are illustrated below the dashed line 207 whereas the server-side processes are illustrated above the dashed line 207. The assistant orchestrator 206 may also determine to respond to the user input by using both the client-side processes and the server-side processes simultaneously. Although FIG. 2 illustrates the assistant orchestrator 206 as being a client-side process, the assistant orchestrator 206 may be a server-side process or may be a hybrid process split between client- and server-side processes.

In particular embodiments, the server-side processes may be as follows after audio data is generated from the system audio API 202. The assistant orchestrator 206 may send the audio data to a remote computing system that hosts different modules of the assistant system 140 to respond to the user input. In particular embodiments, the audio data may be received at a remote automatic speech recognition (ASR) module 208. The ASR module 208 may allow a user to

dictate and have speech transcribed as written text, have a document synthesized as an audio stream, or issue commands that are recognized as such by the system. The ASR module **208** may use statistical models to determine the most likely sequences of words that correspond to a given portion of speech received by the assistant system **140** as audio input. The models may include one or more of hidden Markov models, neural networks, deep learning models, or any combination thereof. The received audio input may be encoded into digital data at a particular sampling rate (e.g., 16, 44.1, or 96 kHz) and with a particular number of bits representing each sample (e.g., 8, 16, or 24 bits).

In particular embodiments, the ASR module **208** may comprise different components. The ASR module **208** may comprise one or more of a grapheme-to-phoneme (G2P) model, a pronunciation learning model, a personalized acoustic model, a personalized language model (PLM), or an end-pointing model. In particular embodiments, the G2P model may be used to determine a user's grapheme-to-phoneme style, e.g., what it may sound like when a particular user speaks a particular word. The personalized acoustic model may be a model of the relationship between audio signals and the sounds of phonetic units in the language. Therefore, such personalized acoustic model may identify how a user's voice sounds. The personalized acoustical model may be generated using training data such as training speech received as audio input and the corresponding phonetic units that correspond to the speech. The personalized acoustical model may be trained or refined using the voice of a particular user to recognize that user's speech. In particular embodiments, the personalized language model may then determine the most likely phrase that corresponds to the identified phonetic units for a particular audio input. The personalized language model may be a model of the probabilities that various word sequences may occur in the language. The sounds of the phonetic units in the audio input may be matched with word sequences using the personalized language model, and greater weights may be assigned to the word sequences that are more likely to be phrases in the language. The word sequence having the highest weight may be then selected as the text that corresponds to the audio input. In particular embodiments, the personalized language model may be also used to predict what words a user is most likely to say given a context. In particular embodiments, the end-pointing model may detect when the end of an utterance is reached.

In particular embodiments, the output of the ASR module **208** may be sent to a remote natural-language understanding (NLU) module **210**. The NLU module **210** may perform named entity resolution (NER). The NLU module **210** may additionally consider contextual information when analyzing the user input. In particular embodiments, an intent and/or a slot may be an output of the NLU module **210**. An intent may be an element in a pre-defined taxonomy of semantic intentions, which may indicate a purpose of a user interacting with the assistant system **140**. The NLU module **210** may classify a user input into a member of the pre-defined taxonomy, e.g., for the input "Play Beethoven's 5th," the NLU module **210** may classify the input as having the intent [IN:play_music]. In particular embodiments, a domain may denote a social context of interaction, e.g., education, or a namespace for a set of intents, e.g., music. A slot may be a named sub-string corresponding to a character string within the user input, representing a basic semantic entity. For example, a slot for "pizza" may be [SL:dish]. In particular embodiments, a set of valid or expected named slots may be conditioned on the classified intent. As an

example and not by way of limitation, for the intent [IN:play_music], a valid slot may be [SL:song_name]. In particular embodiments, the NLU module **210** may additionally extract information from one or more of a social graph, a knowledge graph, or a concept graph, and retrieve a user's profile from one or more remote data stores **212**. The NLU module **210** may further process information from these different sources by determining what information to aggregate, annotating n-grams of the user input, ranking the n-grams with confidence scores based on the aggregated information, and formulating the ranked n-grams into features that can be used by the NLU module **210** for understanding the user input.

In particular embodiments, the NLU module **210** may identify one or more of a domain, an intent, or a slot from the user input in a personalized and context-aware manner. As an example and not by way of limitation, a user input may comprise "show me how to get to the coffee shop". The NLU module **210** may identify the particular coffee shop that the user wants to go based on the user's personal information and the associated contextual information. In particular embodiments, the NLU module **210** may comprise a lexicon of a particular language and a parser and grammar rules to partition sentences into an internal representation. The NLU module **210** may also comprise one or more programs that perform naive semantics or stochastic semantic analysis to the use of pragmatics to understand a user input. In particular embodiments, the parser may be based on a deep learning architecture comprising multiple long-short term memory (LSTM) networks. As an example and not by way of limitation, the parser may be based on a recurrent neural network grammar (RNNG) model, which is a type of recurrent and recursive LSTM algorithm. More information on natural-language understanding may be found in U.S. patent application Ser. No. 16/011,062, filed 18 Jun. 2018, U.S. patent application Ser. No. 16/025,317, filed 2 Jul. 2018, and U.S. patent application Ser. No. 16/038,120, filed 17 Jul. 2018, each of which is incorporated by reference.

In particular embodiments, the output of the NLU module **210** may be sent to a remote reasoning module **214**. The reasoning module **214** may comprise a dialog manager and an entity resolution component. In particular embodiments, the dialog manager may have complex dialog logic and product-related business logic. The dialog manager may manage the dialog state and flow of the conversation between the user and the assistant system **140**. The dialog manager may additionally store previous conversations between the user and the assistant system **140**. In particular embodiments, the dialog manager may communicate with the entity resolution component to resolve entities associated with the one or more slots, which supports the dialog manager to advance the flow of the conversation between the user and the assistant system **140**. In particular embodiments, the entity resolution component may access one or more of the social graph, the knowledge graph, or the concept graph when resolving the entities. Entities may include, for example, unique users or concepts, each of which may have a unique identifier (ID). As an example and not by way of limitation, the knowledge graph may comprise a plurality of entities. Each entity may comprise a single record associated with one or more attribute values. The particular record may be associated with a unique entity identifier. Each record may have diverse values for an attribute of the entity. Each attribute value may be associated with a confidence probability. A confidence probability for an attribute value represents a probability that the value is

accurate for the given attribute. Each attribute value may be also associated with a semantic weight. A semantic weight for an attribute value may represent how the value semantically appropriate for the given attribute considering all the available information. For example, the knowledge graph may comprise an entity of a book “Alice’s Adventures”, which includes information that has been extracted from multiple content sources (e.g., an online social network, online encyclopedias, book review sources, media databases, and entertainment content sources), and then deduped, resolved, and fused to generate the single unique record for the knowledge graph. The entity may be associated with a “fantasy” attribute value which indicates the genre of the book “Alice’s Adventures”. More information on the knowledge graph may be found in U.S. patent application Ser. No. 16/048,049, filed 27 Jul. 2018, and U.S. patent application Ser. No. 16/048,101, filed 27 Jul. 2018, each of which is incorporated by reference.

In particular embodiments, the entity resolution component may check the privacy constraints to guarantee that the resolving of the entities does not violate privacy policies. As an example and not by way of limitation, an entity to be resolved may be another user who specifies in his/her privacy settings that his/her identity should not be searchable on the online social network, and thus the entity resolution component may not return that user’s identifier in response to a request. Based on the information obtained from the social graph, the knowledge graph, the concept graph, and the user profile, and subject to applicable privacy policies, the entity resolution component may therefore resolve the entities associated with the user input in a personalized, context-aware, and privacy-aware manner. In particular embodiments, each of the resolved entities may be associated with one or more identifiers hosted by the social-networking system **160**. As an example and not by way of limitation, an identifier may comprise a unique user identifier (ID) corresponding to a particular user (e.g., a unique username or user ID number). In particular embodiments, each of the resolved entities may be also associated with a confidence score. More information on resolving entities may be found in U.S. patent application Ser. No. 16/048,049, filed 27 Jul. 2018, and U.S. patent application Ser. No. 16/048,072, filed 27 Jul. 2018, each of which is incorporated by reference.

In particular embodiments, the dialog manager may conduct dialog optimization and assistant state tracking. Dialog optimization is the problem of using data to understand what the most likely branching in a dialog should be. As an example and not by way of limitation, with dialog optimization the assistant system **140** may not need to confirm who a user wants to call because the assistant system **140** has high confidence that a person inferred based on dialog optimization would be very likely whom the user wants to call. In particular embodiments, the dialog manager may use reinforcement learning for dialog optimization. Assistant state tracking aims to keep track of a state that changes over time as a user interacts with the world and the assistant system **140** interacts with the user. As an example and not by way of limitation, assistant state tracking may track what a user is talking about, whom the user is with, where the user is, what tasks are currently in progress, and where the user’s gaze is at, etc., subject to applicable privacy policies. In particular embodiments, the dialog manager may use a set of operators to track the dialog state. The operators may comprise the necessary data and logic to update the dialog state. Each operator may act as delta of the dialog state after processing an incoming request. In particular embodiments,

the dialog manager may further comprise a dialog state tracker and an action selector. In alternative embodiments, the dialog state tracker may replace the entity resolution component and resolve the references/mentions and keep track of the state.

In particular embodiments, the reasoning module **214** may further conduct false trigger mitigation. The goal of false trigger mitigation is to detect false triggers (e.g., wake-word) of assistance requests and to avoid generating false records when a user actually does not intend to invoke the assistant system **140**. As an example and not by way of limitation, the reasoning module **214** may achieve false trigger mitigation based on a nonsense detector. If the nonsense detector determines that a wake-word makes no sense at this point in the interaction with the user, the reasoning module **214** may determine that inferring the user intended to invoke the assistant system **140** may be incorrect. In particular embodiments, the output of the reasoning module **214** may be sent a remote dialog arbitrator **216**.

In particular embodiments, each of the ASR module **208**, NLU module **210**, and reasoning module **214** may access the remote data store **212**, which comprises user episodic memories to determine how to assist a user more effectively. More information on episodic memories may be found in U.S. patent application Ser. No. 16/552,559, filed 27 Aug. 2019, which is incorporated by reference. The data store **212** may additionally store the user profile of the user. The user profile of the user may comprise user profile data including demographic information, social information, and contextual information associated with the user. The user profile data may also include user interests and preferences on a plurality of topics, aggregated through conversations on news feed, search logs, messaging platforms, etc. The usage of a user profile may be subject to privacy constraints to ensure that a user’s information can be used only for his/her benefit, and not shared with anyone else. More information on user profiles may be found in U.S. patent application Ser. No. 15/967,239, filed 30 Apr. 2018, which is incorporated by reference.

In particular embodiments, parallel to the aforementioned server-side process involving the ASR module **208**, NLU module **210**, and reasoning module **214**, the client-side process may be as follows. In particular embodiments, the output of the assistant orchestrator **206** may be sent to a local ASR module **216** on the client system **130**. The ASR module **216** may comprise a personalized language model (PLM), a G2P model, and an end-pointing model. Because of the limited computing power of the client system **130**, the assistant system **140** may optimize the personalized language model at run time during the client-side process. As an example and not by way of limitation, the assistant system **140** may pre-compute a plurality of personalized language models for a plurality of possible subjects a user may talk about. When a user requests assistance, the assistant system **140** may then swap these pre-computed language models quickly so that the personalized language model may be optimized locally by the assistant system **140** at run time based on user activities. As a result, the assistant system **140** may have a technical advantage of saving computational resources while efficiently determining what the user may be talking about. In particular embodiments, the assistant system **140** may also re-learn user pronunciations quickly at run time.

In particular embodiments, the output of the ASR module **216** may be sent to a local NLU module **218**. In particular embodiments, the NLU module **218** herein may be more compact compared to the remote NLU module **210** sup-

ported on the server-side. When the ASR module **216** and NLU module **218** process the user input, they may access a local assistant memory **220**. The local assistant memory **220** may be different from the user memories stored on the data store **212** for the purpose of protecting user privacy. In particular embodiments, the local assistant memory **220** may be syncing with the user memories stored on the data store **212** via the network **110**. As an example and not by way of limitation, the local assistant memory **220** may sync a calendar on a user's client system **130** with a server-side calendar associate with the user. In particular embodiments, any secured data in the local assistant memory **220** may be only accessible to the modules of the assistant system **140** that are locally executing on the client system **130**.

In particular embodiments, the output of the NLU module **218** may be sent to a local reasoning module **222**. The reasoning module **222** may comprise a dialog manager and an entity resolution component. Due to the limited computing power, the reasoning module **222** may conduct on-device learning that is based on learning algorithms particularly tailored for client systems **130**. As an example and not by way of limitation, federated learning may be used by the reasoning module **222**. Federated learning is a specific category of distributed machine learning approaches which trains machine learning models using decentralized data residing on end devices such as mobile phones. In particular embodiments, the reasoning module **222** may use a particular federated learning model, namely federated user representation learning, to extend existing neural-network personalization techniques to federated learning. Federated user representation learning can personalize models in federated learning by learning task-specific user representations (i.e., embeddings) or by personalizing model weights. Federated user representation learning is a simple, scalable, privacy-preserving, and resource-efficient. Federated user representation learning may divide model parameters into federated and private parameters. Private parameters, such as private user embeddings, may be trained locally on a client system **130** instead of being transferred to or averaged on a remote server. Federated parameters, by contrast, may be trained remotely on the server. In particular embodiments, the reasoning module **222** may use another particular federated learning model, namely active federated learning to transmit a global model trained on the remote server to client systems **130** and calculate gradients locally on these client systems **130**. Active federated learning may enable the reasoning module to minimize the transmission costs associated with downloading models and uploading gradients. For active federated learning, in each round client systems are selected not uniformly at random, but with a probability conditioned on the current model and the data on the client systems to maximize efficiency. In particular embodiments, the reasoning module **222** may use another particular federated learning model, namely federated Adam. Conventional federated learning model may use stochastic gradient descent (SGD) optimizers. By contrast, the federated Adam model may use moment-based optimizers. Instead of using the averaged model directly as what conventional work does, federated Adam model may use the averaged model to compute approximate gradients. These gradients may be then fed into the federated Adam model, which may de-noise stochastic gradients and use a per-parameter adaptive learning rate. Gradients produced by federated learning may be even noisier than stochastic gradient descent (because data may be not independent and identically distributed), so federated Adam model may help even more deal with the noise. The federated Adam model may use the gradients to take smarter

steps towards minimizing the objective function. The experiments show that conventional federated learning on a benchmark has 1.6% drop in ROC (Receiver Operating Characteristics) curve whereas federated Adam model has only 0.4% drop. In addition, federated Adam model has no increase in communication or on-device computation. In particular embodiments, the reasoning module **222** may also perform false trigger mitigation. This false trigger mitigation may help detect false activation requests, e.g., wake-word, on the client system **130** when the user's speech input comprises data that is subject to privacy constraints. As an example and not by way of limitation, when a user is in a voice call, the user's conversation is private and the false trigger detection based on such conversation can only occur locally on the user's client system **130**.

In particular embodiments, the assistant system **140** may comprise a local context engine **224**. The context engine **224** may process all the other available signals to provide more informative cues to the reasoning module **222**. As an example and not by way of limitation, the context engine **224** may have information related to people, sensory data from client system **130** sensors (e.g., microphone, camera) that are further analyzed by computer vision technologies, geometry constructions, activity data, inertial data (e.g., collected by a VR headset), location, etc. In particular embodiments, the computer vision technologies may comprise human skeleton reconstruction, face detection, facial recognition, hand tracking, eye tracking, etc. In particular embodiments, geometry constructions may comprise constructing objects surrounding a user using data collected by a client system **130**. As an example and not by way of limitation, the user may be wearing AR glasses and geometry construction may aim to determine where the floor is, where the wall is, where the user's hands are, etc. In particular embodiments, inertial data may be data associated with linear and angular motions. As an example and not by way of limitation, inertial data may be captured by AR glasses which measures how a user's body parts move.

In particular embodiments, the output of the local reasoning module **222** may be sent to the dialog arbitrator **216**. The dialog arbitrator **216** may function differently in three scenarios. In the first scenario, the assistant orchestrator **206** determines to use server-side process, for which the dialog arbitrator **216** may transmit the output of the reasoning module **214** to a remote action execution module **226**. In the second scenario, the assistant orchestrator **206** determines to use both server-side processes and client-side processes, for which the dialog arbitrator **216** may aggregate output from both reasoning modules (i.e., remote reasoning module **214** and local reasoning module **222**) of both processes and analyze them. As an example and not by way of limitation, the dialog arbitrator **216** may perform ranking and select the best reasoning result for responding to the user input. In particular embodiments, the dialog arbitrator **216** may further determine whether to use agents on the server-side or on the client-side to execute relevant tasks based on the analysis. In the third scenario, the assistant orchestrator **206** determines to use client-side processes and the dialog arbitrator **216** needs to evaluate the output of the local reasoning module **222** to determine if the client-side processes can complete the task of handling the user input.

In particular embodiments, for the first and second scenarios mentioned above, the dialog arbitrator **216** may determine that the agents on the server-side are necessary to execute tasks responsive to the user input. Accordingly, the dialog arbitrator **216** may send necessary information regarding the user input to the action execution module **226**.

The action execution module **226** may call one or more agents to execute the tasks. In alternative embodiments, the action selector of the dialog manager may determine actions to execute and instruct the action execution module **226** accordingly. In particular embodiments, an agent may be an implementation that serves as a broker across a plurality of content providers for one domain. A content provider may be an entity responsible for carrying out an action associated with an intent or completing a task associated with the intent. In particular embodiments, the agents may comprise first-party agents and third-party agents. In particular embodiments, first-party agents may comprise internal agents that are accessible and controllable by the assistant system **140** (e.g. agents associated with services provided by the online social network, such as messaging services or photo-share services). In particular embodiments, third-party agents may comprise external agents that the assistant system **140** has no control over (e.g., third-party online music application agents, ticket sales agents). The first-party agents may be associated with first-party providers that provide content objects and/or services hosted by the social-networking system **160**. The third-party agents may be associated with third-party providers that provide content objects and/or services hosted by the third-party system **170**. In particular embodiments, each of the first-party agents or third-party agents may be designated for a particular domain. As an example and not by way of limitation, the domain may comprise weather, transportation, music, etc. In particular embodiments, the assistant system **140** may use a plurality of agents collaboratively to respond to a user input. As an example and not by way of limitation, the user input may comprise “direct me to my next meeting.” The assistant system **140** may use a calendar agent to retrieve the location of the next meeting. The assistant system **140** may then use a navigation agent to direct the user to the next meeting.

In particular embodiments, for the second and third scenarios mentioned above, the dialog arbitrator **216** may determine that the agents on the client-side are capable of executing tasks responsive to the user input but additional information is needed (e.g., response templates) or that the tasks can be only handled by the agents on the server-side. If the dialog arbitrator **216** determines that the tasks can be only handled by the agents on the server-side, the dialog arbitrator **216** may send necessary information regarding the user input to the action execution module **226**. If the dialog arbitrator **216** determines that the agents on the client-side are capable of executing tasks but response templates are needed, the dialog arbitrator **216** may send necessary information regarding the user input to a remote response template generation module **228**. The output of the response template generation module **228** may be further sent to a local action execution module **230** executing on the client system **130**.

In particular embodiments, the action execution module **230** may call local agents to execute tasks. A local agent on the client system **130** may be able to execute simpler tasks compared to an agent on the server-side. As an example and not by way of limitation, multiple device-specific implementations (e.g., real-time calls for a client system **130** or a messaging application on the client system **130**) may be handled internally by a single agent. Alternatively, these device-specific implementations may be handled by multiple agents associated with multiple domains. In particular embodiments, the action execution module **230** may additionally perform a set of general executable dialog actions. The set of executable dialog actions may interact with agents, users and the assistant system **140** itself. These

dialog actions may comprise dialog actions for slot request, confirmation, disambiguation, agent execution, etc. The dialog actions may be independent of the underlying implementation of the action selector or dialog policy. Both tree-based policy and model-based policy may generate the same basic dialog actions, with a callback function hiding any action selector specific implementation details.

In particular embodiments, the output from the remote action execution module **226** on the server-side may be sent to a remote response execution module **232**. In particular embodiments, the action execution module **226** may communicate back to the dialog arbitrator **216** for more information. The response execution module **232** may be based on a remote conversational understanding (CU) composer. In particular embodiments, the output from the action execution module **226** may be formulated as a <k, c, u, d> tuple, in which k indicates a knowledge source, c indicates a communicative goal, u indicates a user model, and d indicates a discourse model. In particular embodiments, the CU composer may comprise a natural-language generation (NLG) module and a user interface (UI) payload generator. The natural-language generator may generate a communication content based on the output of the action execution module **226** using different language models and/or language templates. In particular embodiments, the generation of the communication content may be application specific and also personalized for each user. The CU composer may also determine a modality of the generated communication content using the UI payload generator. In particular embodiments, the NLG module may comprise a content determination component, a sentence planner, and a surface realization component. The content determination component may determine the communication content based on the knowledge source, communicative goal, and the user’s expectations. As an example and not by way of limitation, the determining may be based on a description logic. The description logic may comprise, for example, three fundamental notions which are individuals (representing objects in the domain), concepts (describing sets of individuals), and roles (representing binary relations between individuals or concepts). The description logic may be characterized by a set of constructors that allow the natural-language generator to build complex concepts/roles from atomic ones. In particular embodiments, the content determination component may perform the following tasks to determine the communication content. The first task may comprise a translation task, in which the input to the natural-language generator may be translated to concepts. The second task may comprise a selection task, in which relevant concepts may be selected among those resulted from the translation task based on the user model. The third task may comprise a verification task, in which the coherence of the selected concepts may be verified. The fourth task may comprise an instantiation task, in which the verified concepts may be instantiated as an executable file that can be processed by the natural-language generator. The sentence planner may determine the organization of the communication content to make it human understandable. The surface realization component may determine specific words to use, the sequence of the sentences, and the style of the communication content. The UI payload generator may determine a preferred modality of the communication content to be presented to the user. In particular embodiments, the CU composer may check privacy constraints associated with the user to make sure the generation of the communication content follows the privacy policies. More information on natural-language generation may be found in U.S. patent application Ser. No.

15/967,279, filed 30 Apr. 2018, and U.S. patent application Ser. No. 15/966,455, filed 30 Apr. 2018, each of which is incorporated by reference.

In particular embodiments, the output from the local action execution module **230** on the client system **130** may be sent to a local response execution module **234**. The response execution module **234** may be based on a local conversational understanding (CU) composer. The CU composer may comprise a natural-language generation (NLG) module. As the computing power of a client system **130** may be limited, the NLG module may be simple for the consideration of computational efficiency. Because the NLG module may be simple, the output of the response execution module **234** may be sent to a local response expansion module **236**. The response expansion module **236** may further expand the result of the response execution module **234** to make a response more natural and contain richer semantic information.

In particular embodiments, if the user input is based on audio signals, the output of the response execution module **232** on the server-side may be sent to a remote text-to-speech (TTS) module **238**. Similarly, the output of the response expansion module **236** on the client-side may be sent to a local TTS module **240**. Both TTS modules may convert a response to audio signals. In particular embodiments, the output from the response execution module **232**, the response expansion module **236**, or the TTS modules on both sides, may be finally sent to a local render output module **242**. The render output module **242** may generate a response that is suitable for the client system **130**. As an example and not by way of limitation, the output of the response execution module **232** or the response expansion module **236** may comprise one or more of natural-language strings, speech, actions with parameters, or rendered images or videos that can be displayed in a VR headset or AR smart glasses. As a result, the render output module **242** may determine what tasks to perform based on the output of CU composer to render the response appropriately for displaying on the VR headset or AR smart glasses. For example, the response may be visual-based modality (e.g., an image or a video clip) that can be displayed via the VR headset or AR smart glasses. As another example, the response may be audio signals that can be played by the user via VR headset or AR smart glasses. As yet another example, the response may be augmented-reality data that can be rendered VR headset or AR smart glasses for enhancing user experience.

In particular embodiments, the assistant system **140** may have a variety of capabilities including audio cognition, visual cognition, signals intelligence, reasoning, and memories. In particular embodiments, the capability of audio recognition may enable the assistant system **140** to understand a user's input associated with various domains in different languages, understand a conversation and be able to summarize it, perform on-device audio cognition for complex commands, identify a user by voice, extract topics from a conversation and auto-tag sections of the conversation, enable audio interaction without a wake-word, filter and amplify user voice from ambient noise and conversations, understand which client system **130** (if multiple client systems **130** are in vicinity) a user is talking to.

In particular embodiments, the capability of visual cognition may enable the assistant system **140** to perform face detection and tracking, recognize a user, recognize most people of interest in major metropolitan areas at varying angles, recognize majority of interesting objects in the world through a combination of existing machine-learning models and one-shot learning, recognize an interesting moment and

auto-capture it, achieve semantic understanding over multiple visual frames across different episodes of time, provide platform support for additional capabilities in people, places, objects recognition, recognize full set of settings and micro-locations including personalized locations, recognize complex activities, recognize complex gestures to control a client system **130**, handle images/videos from egocentric cameras (e.g., with motion, capture angles, resolution, etc.), accomplish similar level of accuracy and speed regarding images with lower resolution, conduct one-shot registration and recognition of people, places, and objects, and perform visual recognition on a client system **130**.

In particular embodiments, the assistant system **140** may leverage computer vision techniques to achieve visual cognition. Besides computer vision techniques, the assistant system **140** may explore options that can supplement these techniques to scale up the recognition of objects. In particular embodiments, the assistant system **140** may use supplemental signals such as optical character recognition (OCR) of an object's labels, GPS signals for places recognition, signals from a user's client system **130** to identify the user. In particular embodiments, the assistant system **140** may perform general scene recognition (home, work, public space, etc.) to set context for the user and reduce the computer-vision search space to identify top likely objects or people. In particular embodiments, the assistant system **140** may guide users to train the assistant system **140**. For example, crowdsourcing may be used to get users to tag and help the assistant system **140** recognize more objects over time. As another example, users can register their personal objects as part of initial setup when using the assistant system **140**. The assistant system **140** may further allow users to provide positive/negative signals for objects they interact with to train and improve personalized models for them.

In particular embodiments, the capability of signals intelligence may enable the assistant system **140** to determine user location, understand date/time, determine family locations, understand users' calendars and future desired locations, integrate richer sound understanding to identify setting/context through sound alone, build signals intelligence models at run time which may be personalized to a user's individual routines.

In particular embodiments, the capability of reasoning may enable the assistant system **140** to have the ability to pick up any previous conversation threads at any point in the future, synthesize all signals to understand micro and personalized context, learn interaction patterns and preferences from users' historical behavior and accurately suggest interactions that they may value, generate highly predictive proactive suggestions based on micro-context understanding, understand what content a user may want to see at what time of a day, understand the changes in a scene and how that may impact the user's desired content.

In particular embodiments, the capabilities of memories may enable the assistant system **140** to remember which social connections a user previously called or interacted with, write into memory and query memory at will (i.e., open dictation and auto tags), extract richer preferences based on prior interactions and long-term learning, remember a user's life history, extract rich information from egocentric streams of data and auto catalog, and write to memory in structured form to form rich short, episodic and long-term memories.

FIG. 3 illustrates an example diagram flow of server-side processes of the assistant system **140**. In particular embodiments, a server-assistant service module **301** may access a

request manager 302 upon receiving a user request. In alternative embodiments, the user request may be first processed by the remote ASR module 208 if the user request is based on audio signals. In particular embodiments, the request manager 302 may comprise a context extractor 303 and a conversational understanding object generator (CU object generator) 304. The context extractor 303 may extract contextual information associated with the user request. The context extractor 303 may also update contextual information based on the assistant application 136 executing on the client system 130. As an example and not by way of limitation, the update of contextual information may comprise content items are displayed on the client system 130. As another example and not by way of limitation, the update of contextual information may comprise whether an alarm is set on the client system 130. As another example and not by way of limitation, the update of contextual information may comprise whether a song is playing on the client system 130. The CU object generator 304 may generate particular content objects relevant to the user request. The content objects may comprise dialog-session data and features associated with the user request, which may be shared with all the modules of the assistant system 140. In particular embodiments, the request manager 302 may store the contextual information and the generated content objects in data store 212 which is a particular data store implemented in the assistant system 140.

In particular embodiments, the request manager 302 may send the generated content objects to the remote NLU module 210. The NLU module 210 may perform a plurality of steps to process the content objects. At step 305, the NLU module 210 may generate a whitelist for the content objects. In particular embodiments, the whitelist may comprise interpretation data matching the user request. At step 306, the NLU module 210 may perform a featurization based on the whitelist. At step 307, the NLU module 210 may perform domain classification/selection on user request based on the features resulted from the featurization to classify the user request into predefined domains. The domain classification/selection results may be further processed based on two related procedures. At step 308a, the NLU module 210 may process the domain classification/selection result using an intent classifier. The intent classifier may determine the user's intent associated with the user request. In particular embodiments, there may be one intent classifier for each domain to determine the most possible intents in a given domain. As an example and not by way of limitation, the intent classifier may be based on a machine-learning model that may take the domain classification/selection result as input and calculate a probability of the input being associated with a particular predefined intent. At step 308b, the NLU module 210 may process the domain classification/selection result using a meta-intent classifier. The meta-intent classifier may determine categories that describe the user's intent. In particular embodiments, intents that are common to multiple domains may be processed by the meta-intent classifier. As an example and not by way of limitation, the meta-intent classifier may be based on a machine-learning model that may take the domain classification/selection result as input and calculate a probability of the input being associated with a particular predefined meta-intent. At step 309a, the NLU module 210 may use a slot tagger to annotate one or more slots associated with the user request. In particular embodiments, the slot tagger may annotate the one or more slots for the n-grams of the user request. At step 309b, the NLU module 210 may use a meta slot tagger to annotate one or more slots for the classification

result from the meta-intent classifier. In particular embodiments, the meta slot tagger may tag generic slots such as references to items (e.g., the first), the type of slot, the value of the slot, etc. As an example and not by way of limitation, a user request may comprise "change 500 dollars in my account to Japanese yen." The intent classifier may take the user request as input and formulate it into a vector. The intent classifier may then calculate probabilities of the user request being associated with different predefined intents based on a vector comparison between the vector representing the user request and the vectors representing different predefined intents. In a similar manner, the slot tagger may take the user request as input and formulate each word into a vector. The intent classifier may then calculate probabilities of each word being associated with different predefined slots based on a vector comparison between the vector representing the word and the vectors representing different predefined slots. The intent of the user may be classified as "changing money". The slots of the user request may comprise "500", "dollars", "account", and "Japanese yen". The meta-intent of the user may be classified as "financial service". The meta slot may comprise "finance".

In particular embodiments, the NLU module 210 may comprise a semantic information aggregator 310. The semantic information aggregator 310 may help the NLU module 210 improve the domain classification/selection of the content objects by providing semantic information. In particular embodiments, the semantic information aggregator 310 may aggregate semantic information in the following way. The semantic information aggregator 310 may first retrieve information from a user context engine 315. In particular embodiments, the user context engine 315 may comprise offline aggregators and an online inference service. The offline aggregators may process a plurality of data associated with the user that are collected from a prior time window. As an example and not by way of limitation, the data may include news feed posts/comments, interactions with news feed posts/comments, search history, etc., that are collected during a predetermined timeframe (e.g., from a prior 90-day window). The processing result may be stored in the user context engine 315 as part of the user profile. The online inference service may analyze the conversational data associated with the user that are received by the assistant system 140 at a current time. The analysis result may be stored in the user context engine 315 also as part of the user profile. In particular embodiments, both the offline aggregators and online inference service may extract personalization features from the plurality of data. The extracted personalization features may be used by other modules of the assistant system 140 to better understand user input. In particular embodiments, the semantic information aggregator 310 may then process the retrieved information, i.e., a user profile, from the user context engine 315 in the following steps. At step 311, the semantic information aggregator 310 may process the retrieved information from the user context engine 315 based on natural-language processing (NLP). In particular embodiments, the semantic information aggregator 310 may tokenize text by text normalization, extract syntax features from text, and extract semantic features from text based on NLP. The semantic information aggregator 310 may additionally extract features from contextual information, which is accessed from dialog history between a user and the assistant system 140. The semantic information aggregator 310 may further conduct global word embedding, domain-specific embedding, and/or dynamic embedding based on the contextual information. At step 312, the processing result may be annotated

with entities by an entity tagger. Based on the annotations, the semantic information aggregator **310** may generate dictionaries for the retrieved information at step **313**. In particular embodiments, the dictionaries may comprise global dictionary features which can be updated dynamically 5 offline. At step **314**, the semantic information aggregator **310** may rank the entities tagged by the entity tagger. In particular embodiments, the semantic information aggregator **310** may communicate with different graphs **320** including one or more of the social graph, the knowledge graph, or the concept graph to extract ontology data that is relevant to the retrieved information from the user context engine **315**. In particular embodiments, the semantic information aggregator **310** may aggregate the user profile, the ranked entities, and the information from the graphs **320**. The semantic information aggregator **310** may then provide the aggregated information to the NLU module **210** to facilitate the domain classification/selection.

In particular embodiments, the output of the NLU module **210** may be sent to the remote reasoning module **214**. The reasoning module **214** may comprise a co-reference component **325**, an entity resolution component **330**, and a dialog manager **335**. The output of the NLU module **210** may be first received at the co-reference component **325** to interpret references of the content objects associated with the user request. In particular embodiments, the co-reference component **325** may be used to identify an item to which the user request refers. The co-reference component **325** may comprise reference creation **326** and reference resolution **327**. In particular embodiments, the reference creation **326** may create references for entities determined by the NLU module **210**. The reference resolution **327** may resolve these references accurately. As an example and not by way of limitation, a user request may comprise “find me the nearest grocery store and direct me there”. The co-reference component **325** may interpret “there” as “the nearest grocery store”. In particular embodiments, the co-reference component **325** may access the user context engine **315** and the dialog manager **335** when necessary to interpret references with improved accuracy.

In particular embodiments, the identified domains, intents, meta-intents, slots, and meta slots, along with the resolved references may be sent to the entity resolution component **330** to resolve relevant entities. The entity resolution component **330** may execute generic and domain-specific entity resolution. In particular embodiments, the entity resolution component **330** may comprise domain entity resolution **331** and generic entity resolution **332**. The domain entity resolution **331** may resolve the entities by categorizing the slots and meta slots into different domains. In particular embodiments, entities may be resolved based on the ontology data extracted from the graphs **320**. The ontology data may comprise the structural relationship between different slots/meta-slots and domains. The ontology may also comprise information of how the slots/meta-slots may be grouped, related within a hierarchy where the higher level comprises the domain, and subdivided according to similarities and differences. The generic entity resolution **332** may resolve the entities by categorizing the slots and meta slots into different generic topics. In particular embodiments, the resolving may be also based on the ontology data extracted from the graphs **320**. The ontology data may comprise the structural relationship between different slots/meta-slots and generic topics. The ontology may also comprise information of how the slots/meta-slots may be grouped, related within a hierarchy where the higher level comprises the topic, and subdivided according to similarities

and differences. As an example and not by way of limitation, in response to the input of an inquiry of the advantages of a particular brand of electric car, the generic entity resolution **332** may resolve the referenced brand of electric car as vehicle and the domain entity resolution **331** may resolve the referenced brand of electric car as electric car.

In particular embodiments, the output of the entity resolution component **330** may be sent to the dialog manager **335** to advance the flow of the conversation with the user. The dialog manager **335** may be an asynchronous state machine that repeatedly updates the state and selects actions based on the new state. The dialog manager **335** may comprise dialog intent resolution **336** and dialog state tracker **337**. In particular embodiments, the dialog manager **335** may execute the selected actions and then call the dialog state tracker **337** again until the action selected requires a user response, or there are no more actions to execute. Each action selected may depend on the execution result from previous actions. In particular embodiments, the dialog intent resolution **336** may resolve the user intent associated with the current dialog session based on dialog history between the user and the assistant system **140**. The dialog intent resolution **336** may map intents determined by the NLU module **210** to different dialog intents. The dialog intent resolution **336** may further rank dialog intents based on signals from the NLU module **210**, the entity resolution component **330**, and dialog history between the user and the assistant system **140**. In particular embodiments, instead of directly altering the dialog state, the dialog state tracker **337** may be a side-effect free component and generate n-best candidates of dialog state update operators that propose updates to the dialog state. The dialog state tracker **337** may comprise intent resolvers containing logic to handle different types of NLU intent based on the dialog state and generate the operators. In particular embodiments, the logic may be organized by intent handler, such as a disambiguation intent handler to handle the intents when the assistant system **140** asks for disambiguation, a confirmation intent handler that comprises the logic to handle confirmations, etc. Intent resolvers may combine the turn intent together with the dialog state to generate the contextual updates for a conversation with the user. A slot resolution component may then recursively resolve the slots in the update operators with resolution providers including the knowledge graph and domain agents. In particular embodiments, the dialog state tracker **337** may update/rank the dialog state of the current dialog session. As an example and not by way of limitation, the dialog state tracker **337** may update the dialog state as “completed” if the dialog session is over. As another example and not by way of limitation, the dialog state tracker **337** may rank the dialog state based on a priority associated with it.

In particular embodiments, the reasoning module **214** may communicate with the remote action execution module **226** and the dialog arbitrator **216**, respectively. In particular embodiments, the dialog manager **335** of the reasoning module **214** may communicate with a task completion component **340** of the action execution module **226** about the dialog intent and associated content objects. In particular embodiments, the task completion module **340** may rank different dialog hypotheses for different dialog intents. The task completion module **340** may comprise an action selector **341**. In alternative embodiments, the action selector **341** may be comprised in the dialog manager **335**. In particular embodiments, the dialog manager **335** may additionally check against dialog policies **345** comprised in the dialog arbitrator **216** regarding the dialog state. In particular

embodiments, a dialog policy 345 may comprise a data structure that describes an execution plan of an action by an agent 350. The dialog policy 345 may comprise a general policy 346 and task policies 347. In particular embodiments, the general policy 346 may be used for actions that are not specific to individual tasks. The general policy 346 may comprise handling low confidence intents, internal errors, unacceptable user response with retries, skipping or inserting confirmation based on ASR or NLU confidence scores, etc. The general policy 346 may also comprise the logic of ranking dialog state update candidates from the dialog state tracker 337 output and pick the one to update (such as picking the top ranked task intent). In particular embodiments, the assistant system 140 may have a particular interface for the general policy 346, which allows for consolidating scattered cross-domain policy/business-rules, especial those found in the dialog state tracker 337, into a function of the action selector 341. The interface for the general policy 346 may also allow for authoring of self-contained sub-policy units that may be tied to specific situations or clients, e.g., policy functions that may be easily switched on or off based on clients, situation, etc. The interface for the general policy 346 may also allow for providing a layering of policies with back-off, i.e. multiple policy units, with highly specialized policy units that deal with specific situations being backed up by more general policies 346 that apply in wider circumstances. In this context the general policy 346 may alternatively comprise intent or task specific policy. In particular embodiments, a task policy 347 may comprise the logic for action selector 341 based on the task and current state. In particular embodiments, there may be the following four types of task policies 347: 1) manually crafted tree-based dialog plans; 2) coded policy that directly implements the interface for generating actions; 3) configurator-specified slot-filling tasks; and 4) machine-learning model based policy learned from data. In particular embodiments, the assistant system 140 may bootstrap new domains with rule-based logic and later refine the task policies 347 with machine-learning models. In particular embodiments, a dialog policy 345 may be a tree-based policy, which is a pre-constructed dialog plan. Based on the current dialog state, a dialog policy 345 may choose a node to execute and generate the corresponding actions. As an example and not by way of limitation, the tree-based policy may comprise topic grouping nodes and dialog action (leaf) nodes.

In particular embodiments, the action selector 341 may take candidate operators of dialog state and consult the dialog policy 345 to decide what action should be executed. The assistant system 140 may use a hierarchical dialog policy with general policy 346 handling the cross-domain business logic and task policies 347 handles the task/domain specific logic. In particular embodiments, the general policy 346 may pick one operator from the candidate operators to update the dialog state, followed by the selection of a user facing action by a task policy 347. Once a task is active in the dialog state, the corresponding task policy 347 may be consulted to select right actions. In particular embodiments, both the dialog state tracker 337 and the action selector 341 may not change the dialog state until the selected action is executed. This may allow the assistant system 140 to execute the dialog state tracker 337 and the action selector 341 for processing speculative ASR results and to do n-best ranking with dry runs. In particular embodiments, the action selector 341 may take the dialog state update operators as part of the input to select the dialog action. The execution of the dialog action may generate a set of expectation to

instruct the dialog state tracker 337 to handler future turns. In particular embodiments, an expectation may be used to provide context to the dialog state tracker 337 when handling the user input from next turn. As an example and not by way of limitation, slot request dialog action may have the expectation of proving a value for the requested slot.

In particular embodiments, the dialog manager 335 may support multi-turn compositional resolution of slot mentions. For a compositional parse from the NLU 210, the resolver may recursively resolve the nested slots. The dialog manager 335 may additionally support disambiguation for the nested slots. As an example and not by way of limitation, the user request may be “remind me to call Alex”. The resolver may need to know which Alex to call before creating an actionable reminder to-do entity. The resolver may halt the resolution and set the resolution state when further user clarification is necessary for a particular slot. The general policy 346 may examine the resolution state and create corresponding dialog action for user clarification. In dialog state tracker 337, based on the user request and the last dialog action, the dialog manager may update the nested slot. This capability may allow the assistant system 140 to interact with the user not only to collect missing slot values but also to reduce ambiguity of more complex/ambiguous utterances to complete the task. In particular embodiments, the dialog manager may further support requesting missing slots in a nested intent and multi-intent user requests (e.g., “take this photo and send it to Dad”). In particular embodiments, the dialog manager 335 may support machine-learning models for more robust dialog experience. As an example and not by way of limitation, the dialog state tracker 337 may use neural network based models (or any other suitable machine-learning models) to model belief over task hypotheses. As another example and not by way of limitation, for action selector 341, highest priority policy units may comprise white-list/black-list overrides, which may have to occur by design; middle priority units may comprise machine-learning models designed for action selection; and lower priority units may comprise rule-based fallbacks when the machine-learning models elect not to handle a situation. In particular embodiments, machine-learning model based general policy unit may help the assistant system 140 reduce redundant disambiguation or confirmation steps, thereby reducing the number of turns to execute the user request.

In particular embodiments, the action execution module 226 may call different agents 350 for task execution. An agent 350 may select among registered content providers to complete the action. The data structure may be constructed by the dialog manager 335 based on an intent and one or more slots associated with the intent. A dialog policy 345 may further comprise multiple goals related to each other through logical operators. In particular embodiments, a goal may be an outcome of a portion of the dialog policy and it may be constructed by the dialog manager 335. A goal may be represented by an identifier (e.g., string) with one or more named arguments, which parameterize the goal. As an example and not by way of limitation, a goal with its associated goal argument may be represented as {confirm_artist, args: {artist: “Madonna”}}. In particular embodiments, a dialog policy may be based on a tree-structured representation, in which goals are mapped to leaves of the tree. In particular embodiments, the dialog manager 335 may execute a dialog policy 345 to determine the next action to carry out. The dialog policies 345 may comprise generic policy 346 and domain specific policies 347, both of which may guide how to select the next system action based on the

dialog state. In particular embodiments, the task completion component 340 of the action execution module 226 may communicate with dialog policies 345 comprised in the dialog arbitrator 216 to obtain the guidance of the next system action. In particular embodiments, the action selection component 341 may therefore select an action based on the dialog intent, the associated content objects, and the guidance from dialog policies 345.

In particular embodiments, the output of the action execution module 226 may be sent to the remote response execution module 232. Specifically, the output of the task completion component 340 of the action execution module 226 may be sent to the CU composer 355 of the response execution module 226. In alternative embodiments, the selected action may require one or more agents 350 to be involved. As a result, the task completion module 340 may inform the agents 350 about the selected action. Meanwhile, the dialog manager 335 may receive an instruction to update the dialog state. As an example and not by way of limitation, the update may comprise awaiting agents' 350 response. In particular embodiments, the CU composer 355 may generate a communication content for the user using a natural-language generation (NLG) module 356 based on the output of the task completion module 340. In particular embodiments, the NLG module 356 may use different language models and/or language templates to generate natural language outputs. The generation of natural language outputs may be application specific. The generation of natural language outputs may be also personalized for each user. The CU composer 355 may also determine a modality of the generated communication content using the UI payload generator 357. Since the generated communication content may be considered as a response to the user request, the CU composer 355 may additionally rank the generated communication content using a response ranker 358. As an example and not by way of limitation, the ranking may indicate the priority of the response.

In particular embodiments, the response execution module 232 may perform different tasks based on the output of the CU composer 355. These tasks may include writing (i.e., storing/updating) the dialog state 361 retrieved from data store 212 and generating responses 362. In particular embodiments, the output of CU composer 355 may comprise one or more of natural-language strings, speech, actions with parameters, or rendered images or videos that can be displayed in a VR headset or AR smart glass. As a result, the response execution module 232 may determine what tasks to perform based on the output of CU composer 355. In particular embodiments, the generated response and the communication content may be sent to the local render output module 242 by the response execution module 232. In alternative embodiments, the output of the CU composer 355 may be additionally sent to the remote TTS module 238 if the determined modality of the communication content is audio. The speech generated by the TTS module 238 and the response generated by the response execution module 232 may be then sent to the render output module 242.

FIG. 4 illustrates an example diagram flow of processing a user input by the assistant system 140. As an example and not by way of limitation, the user input may be based on audio signals. In particular embodiments, a mic array 402 of the client system 130 may receive the audio signals (e.g., speech). The audio signals may be transmitted to a process loop 404 in a format of audio frames. In particular embodiments, the process loop 404 may send the audio frames for voice activity detection (VAD) 406 and wake-on-voice (WoV) detection 408. The detection results may be returned

to the process loop 404. If the WoV detection 408 indicates the user wants to invoke the assistant system 140, the audio frames together with the VAD 406 result may be sent to an encode unit 410 to generate encoded audio data. After encoding, the encoded audio data may be sent to an encrypt unit 412 for privacy and security purpose, followed by a link unit 414 and decrypt unit 416. After decryption, the audio data may be sent to a mic driver 418, which may further transmit the audio data to an audio service module 420. In alternative embodiments, the user input may be received at a wireless device (e.g., Bluetooth device) paired with the client system 130. Correspondingly, the audio data may be sent from a wireless-device driver 422 (e.g., Bluetooth driver) to the audio service module 420. In particular embodiments, the audio service module 420 may determine that the user input can be fulfilled by an application executing on the client system 130. Accordingly, the audio service module 420 may send the user input to a real-time communication (RTC) module 424. The RTC module 424 may deliver audio packets to a video or audio communication system (e.g., VOIP or video call). The RTC module 424 may call a relevant application (App) 426 to execute tasks related to the user input.

In particular embodiments, the audio service module 420 may determine that the user is requesting assistance that needs the assistant system 140 to respond. Accordingly, the audio service module 420 may inform the client-assistant service module 426. In particular embodiments, the client-assistant service module 426 may communicate with the assistant orchestrator 206. The assistant orchestrator 206 may determine whether to use client-side processes or server-side processes to respond to the user input. In particular embodiments, the assistant orchestrator 206 may determine to use client-side processes and inform the client-assistant service module 426 about such decision. As a result, the client-assistant service module 426 may call relevant modules to respond to the user input.

In particular embodiments, the client-assistant service module 426 may use the local ASR module 216 to analyze the user input. The ASR module 216 may comprise a grapheme-to-phoneme (G2P) model, a pronunciation learning model, a personalized language model (PLM), an end-pointing model, and a personalized acoustic model. In particular embodiments, the client-assistant service module 426 may further use the local NLU module 218 to understand the user input. The NLU module 218 may comprise a named entity resolution (NER) component and a contextual session-based NLU component. In particular embodiments, the client-assistant service module 426 may use an intent broker 428 to analyze the user's intent. To be accurate about the user's intent, the intent broker 428 may access an entity store 430 comprising entities associated with the user and the world. In alternative embodiments, the user input may be submitted via an application 432 executing on the client system 130. In this case, an input manager 434 may receive the user input and analyze it by an application environment (App Env) module 436. The analysis result may be sent to the application 432 which may further send the analysis result to the ASR module 216 and NLU module 218. In alternative embodiments, the user input may be directly submitted to the client-assistant service module 426 via an assistant application 438 executing on the client system 130. Then the client-assistant service module 426 may perform similar procedures based on modules as aforementioned, i.e., the ASR module 216, the NLU module 218, and the intent broker 428.

In particular embodiments, the assistant orchestrator **206** may determine to user server-side process. Accordingly, the assistant orchestrator **206** may send the user input to one or more computing systems that host different modules of the assistant system **140**. In particular embodiments, a server-assistant service module **301** may receive the user input from the assistant orchestrator **206**. The server-assistant service module **301** may instruct the remote ASR module **208** to analyze the audio data of the user input. The ASR module **208** may comprise a grapheme-to-phoneme (G2P) model, a pronunciation learning model, a personalized language model (PLM), an end-pointing model, and a personalized acoustic model. In particular embodiments, the server-assistant service module **301** may further instruct the remote NLU module **210** to understand the user input. In particular embodiments, the server-assistant service module **301** may call the remote reasoning model **214** to process the output from the ASR module **208** and the NLU module **210**. In particular embodiments, the reasoning model **214** may perform entity resolution and dialog optimization. In particular embodiments, the output of the reasoning model **314** may be sent to the agent **350** for executing one or more relevant tasks.

In particular embodiments, the agent **350** may access an ontology module **440** to accurately understand the result from entity resolution and dialog optimization so that it can execute relevant tasks accurately. The ontology module **440** may provide ontology data associated with a plurality of predefined domains, intents, and slots. The ontology data may also comprise the structural relationship between different slots and domains. The ontology data may further comprise information of how the slots may be grouped, related within a hierarchy where the higher level comprises the domain, and subdivided according to similarities and differences. The ontology data may also comprise information of how the slots may be grouped, related within a hierarchy where the higher level comprises the topic, and subdivided according to similarities and differences. Once the tasks are executed, the agent **350** may return the execution results together with a task completion indication to the reasoning module **214**.

The embodiments disclosed herein may include or be implemented in conjunction with an artificial reality system. Artificial reality is a form of reality that has been adjusted in some manner before presentation to a user, which may include, e.g., a virtual reality (VR), an augmented reality (AR), a mixed reality (MR), a hybrid reality, or some combination and/or derivatives thereof. Artificial reality content may include completely generated content or generated content combined with captured content (e.g., real-world photographs). The artificial reality content may include video, audio, haptic feedback, or some combination thereof, and any of which may be presented in a single channel or in multiple channels (such as stereo video that produces a three-dimensional effect to the viewer). Additionally, in some embodiments, artificial reality may be associated with applications, products, accessories, services, or some combination thereof, that are, e.g., used to create content in an artificial reality and/or used in (e.g., perform activities in) an artificial reality. The artificial reality system that provides the artificial reality content may be implemented on various platforms, including a head-mounted display (HIVID) connected to a host computer system, a standalone HIVID, a mobile device or computing system, or any other hardware platform capable of providing artificial reality content to one or more viewers.

In-Call Experience Enhancement for Assistant Systems

In particular embodiment, an in-call experience enhancement in which the assistant system is persistently active, but on standby during a call (such as a video or audio call) or other communication session (such as a text message thread), is provided. Such a persistently active assistant system may enable a user to invoke it in real-time during the call to execute tasks related to one or more other users on the call. Furthermore, the persistently active assistant system may allow a single communication domain to be used in which the user can communicate with both other people via the call and with the assistant system itself. Current assistant systems typically go dormant during calls, so that a user must pause the call and reawaken the assistant system in order to issue commands. Thus, this single communication domain may greatly improve the user's experience, enabling a more social and natural interaction. The persistent assistant system may utilize an underlying multimodal architecture having separate context and scene understanding engines. The context engine may also be persistent during the call, gathering data for use by other modules in the assistant system that responds to a user query (subject to privacy settings). By contrast, the scene understanding engine may be awakened as needed to receive the data gathered by the context engine and determines a relationship among detected entities. Accordingly, with a video call in particular serving as a social experience backdrop, this persistent assistant system may enable numerous social, utility, communication, and image processing functionalities to be performed.

Although this disclosure describes providing a persistent assistant system and particular social functions in a particular manner, this disclosure contemplates providing a persistent assistant system and any suitable social functions in any suitable manner.

FIG. 5 illustrates an example multimodal architecture of the assistant system **140**. The persistent assistant system **140** may use such an underlying multimodal architecture with separate context and scene understanding engines, as well as various sensors within the user device hosting the assistant system **140**. Sensor/gaze information **501** and passive audio **502** (e.g., background audio picked up by a microphone on the client system **130**) may be gathered (subject to privacy settings), along with a content vector **504** derived from the image **503** (e.g., input from a camera), and input into the context engine **510**. The context engine **510** may generate a context **511**. This context **511** and/or any gesture information **512** may then be input into a non-verbal intent recognizer **513** to determine an intent of the user input. Detection of a wake-word **514** may trigger capture of active audio **515** (e.g., a user's audio input to the assistant system **140**), which may then be input into an ASR module **208** to generate text **517**. This text **517**, and/or any initial text **518** input by a user, may be transmitted to the NLU module **210** to generate semantic information **519**. Sensor/gaze information **501**, passive audio **502**, image information **503**, active audio **515**, and/or text **517** may be input into the scene understanding engine **520**, which may in turn exchange data with the NLU **210**. Output of the scene understanding engine **520** may then be input into the non-verbal intent recognizer **513** and/or into a dialog module **530**. Data from an assistant user memory **550** may also be input into the non-verbal intent recognizer **513** and/or into the dialog module **530**. An assistant state tracker module **531**, which includes a context tracker **532**, a resolver **533**, and a task state tracker **534**, receives these various data items as input. Output from the assistant state tracker **531** then may be input into an action selector **341**, which in turn may include a general policy unit

346 and a task policy unit 347. The output of the dialog module 530 may be sent to ASR 208 to invoke heavy processing, for example, when a request that requires an understanding of the scene is received. As an example and not by way of limitation, a user viewing an image or video stream may request “tell me more about the dog” to the assistant system 140. Segmenting out objects from the image may be resource intensive and unnecessary to answer most requests, so the heavy processing of segmenting objects (for example, to segment out the dog referred to in the request) may only be performed when requested, rather than continuously.

FIG. 6A illustrates an example initial scene 600 viewed during a video call on a first client system 130 of a first user. In particular embodiments, the assistant system 140 may establish a video call between a plurality of client systems 130. Each client system 130 may be associated with one or more users (e.g., participants in the video call). The assistant system 140 may receive a request from the first client system 130 of the first user identifying one or more other users to add to a video call and may assign a call identifier (ID) to the video call. As an example and not by way of limitation, the assistant system 140 may use this call ID in monitoring the video call and context information of the scene 600 and of various client systems 130 participating in it. In particular embodiments, the assistant system 140 may itself be added as a participant in the video call, subject to privacy settings of each of the users of the video call. As an example and not by way of limitation, a first user may request to add the assistant system 140 to the video call, and, if each of the other users permit, the assistant system 140 may be added as a participant to the video call; otherwise, the assistant system 140 may not be added. In particular embodiments, the assistant system 140 may establish a video call customized for business. As an example and not by way of limitation, the first user may communicate with customer service agents via the video call (for example, to show a defective product they received, or to get help with setting up a new product); this assistant system 140 may conceal identifying information of the first user during this video call for privacy. Although this disclosure describes establishing a video call in a particular manner, this disclosure contemplates establishing a video call in any suitable manner.

In particular embodiments, the assistant system 140 may receive, from the first client system 130 from among the plurality of client systems 130, a request from the first user of the first client system 130 to be performed by the persistent assistant system 140 during the video call. The request may be a manual request, a spoken request, a gesture as a request, other suitable input associated with a request, or any combination thereof. The request may reference one or more second users associated with the plurality of client systems 130 in the video call. As an example and not by way of limitation, the first user Alice may be on a video call with the second users Bob, Sarah, and Carol. Alice may then speak to the assistant system 140 and say “Hey Assistant, take a photo of Bob,” referencing the second user Bob on the video call. Although this disclosure describes receiving user requests in a particular manner, this disclosure contemplates receiving user requests in any suitable manner.

In particular embodiments, the assistant system 140 may receive a wake-word 514 that precedes the request and, in response to receiving this wake-word 514, may send instructions to the first client system 130 to mute the video call at the first client system 130. As an example and not by way of limitation, the first user may say “Hey Assistant,” during a video call and, upon detecting the wake-word “Hey Assis-

tant” indicating that the first user is about to make a request of the assistant system 140, the assistant system 140 may maintain user privacy by muting the video call on the first client system 130 so that other users participating in the video call cannot hear the first user’s request. However, in particular embodiments, the first user may explicitly instruct the first client system 130 to mute the video call (e.g., through a spoken request like “Hey Assistant, mute me,” or the selection of a button or icon for muting the call) while communicating with the assistant system 140. Although this disclosure describes detecting a wake-word in a particular manner, this disclosure contemplates detecting a wake-word in any suitable manner.

In particular embodiments, the assistant system 140 may detect a gaze of the first user directed at one or more entities in the video call and may infer the request based on the gaze. Such gaze detection may be made using, for example, eye tracking. As an example and not by way of limitation, the assistant system may detect that the first user repeatedly looks at a clock and infer that he wants to know if a task on his calendar occurs soon. The assistant system 140 may accordingly inform him, without any manual or spoken input on the part of the first user, that the next task on his calendar (such as a meeting) starts in five minutes, and that he should head toward the meeting location while on the video call. As another example and not by way of limitation, the assistant system 140 may detect that the first user repeatedly looks at a second user, and may infer that the first user wants to know the identity of the second user, or wants the smart camera to focus on that second user. The assistant system 140 may thus accordingly inform the first user as to the identity of the second user, and/or track the second user with the smart camera. Although this disclosure describes inferring a request based on gaze estimation in a particular manner, this disclosure contemplates inferring a request in any suitable manner.

In particular embodiments, the request received by the assistant system 140 may reference one or more second users associated with the plurality of client systems 130 in the video call. The request may explicitly refer to a second user by name, or it may imply a second user through a relationship with another second user (e.g., “who is the person to the left of Alice?”). As an example and not by way of limitation, the request may be an instruction to focus the display of the first client system 130 on one or more of the second users, as discussed below in with respect to FIG. 7A. As another example and not by way of limitation, the request may be an instruction to repeat or summarize speech of one or more of the second users. As another example and not by way of limitation, the request may comprise an instruction to perform a virtual activity with respect to one or more of the second users, as discussed below with respect to FIG. 8. As yet another example and not by way of limitation, the request may be an instruction to share a content item with one or more of the second users, as discussed below with respect to FIG. 9. Although this disclosure describes receiving requests referencing users in a particular manner, this disclosure contemplates receiving requests referencing users in any suitable manner.

Certain technical challenges exist in maintaining a quality video call between users. Video calls may lack a feeling of genuine social interaction; providing more social functions that may be performed during the actual video call may thus increase user interaction and satisfaction with the video call. However, one technical challenge to this may include identifying users in the video call that the first user in the video call wants to perform some social function with, as well as

actually understanding the scene and context of the video call in order to more accurately execute the social function. A solution presented by embodiments disclosed herein to address this challenge may thus include continuously gathering context of the video call via the context engine **510** and feeding this gathered information into the scene understanding engine **520**, in order to generate relationship information between people and objects in the scene of the video call. Further, certain embodiments disclosed herein may provide one or more technical advantages. As an example, accurately identifying users and objects in the video call, as well as their context and relationship information (subject to privacy settings), may enable the first user to perform a variety of social functions with respect to entities in the video call, even when the first user communicates those functions ambiguously.

In particular embodiments, upon receiving the request from the first client system **130** referencing one or more second users, the assistant system **140** may determine an intent of the request and one or more user identifiers (user IDs) of the one or more second users referenced by the request. Accurately identifying users and objects in a video call, as well as their context and relationship information, may enable the first user to perform a variety of social functions with respect to entities in the video call, even when the first user communicates those functions ambiguously. As an example and not by way of limitation, the assistant system **140** may receive and execute a request to focus a camera on a particular person. In particular embodiments, the assistant system **140** may determine the one or more user identifiers of the one or more second users referenced by the request through determination of the second users' respective user IDs or through facial recognition of the second users (subject to privacy settings). Upon recognizing a given user, the assistant system **140** may assign them a user ID. As an example and not by way of limitation, both active users currently using client systems **130** and background users viewable in the frame of the video call may be identified. In particular embodiments, the identified users may be modified dynamically as, for example, people enter and leave the frame of the video call, as discussed below with respect to FIGS. **6B** and **7B**. Although this disclosure describes determining intent and user IDs in a particular manner, this disclosure contemplates determining intent and user IDs in any suitable manner.

In particular embodiments, the assistant system **140** may determine that the intent of the request is to modify a characteristic (e.g., an appearance or voice) of one or more second users. In particular embodiments, this intent may be determined as an explicit command to modify this characteristic. As an example and not by way of limitation, the assistant system **140** may thus modify the appearance of the second user having the identified user identifier by adding a mask or special effects to the second user on the display of the first client system **130**, as discussed below with respect to FIG. **8**. Although this disclosure describes determining an intent in a particular manner, this disclosure contemplates determining an intent in any suitable manner.

In particular embodiments, the assistant system **140** may instruct the assistant system to execute the request based on the determined intent and user IDs. This request may be executed on either or both of the client-side process or the server-side process of a hybrid assistant system. As an example and not by way of limitation, if the intent of the request indicates that additional user information is needed to execute the request, the assistant system **140** may, subject to privacy settings, retrieve user profile information of one

or more of the identified second users in response to this determined intent and the one or more user identifiers and generate the response based on this retrieved user profile information. As an example and not by way of limitation, the user profile information may indicate information of an interest or recent activity of one or more of the second users. To execute the request, the assistant system **140** may make use of the context engine and scene understanding engine of the underlying multimodal architecture. Although this disclosure describes executing requests in a particular manner, this disclosure contemplates executing requests in any suitable manner.

In particular embodiments, the assistant system **140** may access, from the context engine **510** of the assistant system **140**, context data associated with the video call. This context data may indicate properties of a scene of the video call. As an example and not by way of limitation, the context data may indicate identifications of objects within the scene. As another example and not by way of limitation, the context data may indicate user IDs of users within the scene. In particular embodiments, the context engine **510** may analyze these properties of the scene in real time during the actual video call. The context engine **510** may analyze these properties through facial, activity, or object recognition, and enter the detected context data into chart **650**. As an example and not by way of limitation, the assistant system **140** may determine the identities of particular people (Alice and Bob), their activities (standing and speaking), and their location. Although this disclosure describes accessing context information in a particular manner, this disclosure contemplates accessing context information in any suitable manner.

In particular embodiments, the context engine **510** may always be on during the video call, gathering intelligence for use later in the pipeline of multimodal architecture **500** that responds to a user query. The context engine **510** may thus function as a sort of ambient mode of the assistant system **140**, constantly monitoring the video call as well as the first user and capturing information that may be needed to respond to a future user request. Using sensors such as a smart camera, which may also be always on, the context engine **510** may identify particular objects (e.g., a water bottle in user Alice's hand), activities (e.g., whether a user is cooking or standing), or locations (e.g., whether the user is at a museum or a concert) in scene **600** viewed during the video call. The context engine **510** may further continuously track the first user of client device **130** himself, such as through eye tracking/gaze estimation. Although this disclosure describes operating a context engine in a particular manner, this disclosure contemplates operating a context engine in any suitable manner.

FIG. **6B** illustrates an example chart **650** of information of the scene **600** generated by the always-on context engine **510**. This chart **650** may be the output of context engine **510** when analyzing scene **600**. The chart **650** may include various categories, such as social presence **651**, user activity class **652**, focal object recognition **653**, user location **654**, or significant events detection **655**. The social presence category **651** may include social information of people in the scene **600** of the video call, allowing particular individuals to be recognized. As some functionalities of the assistant system **140** may require it to be able identify particular people within the scene **600** (for example, informing the first user of the name of an unknown second user in the scene or modifying the appearance of a particular second user), user recognition may enable these user-specific functions to be performed. User activity class **652** may indicate current activity of a detected user, classified into a taxonomy of

activity classes; and user location **654** may indicate deeper knowledge information about the location of a user on a personal, group, or world-knowledge basis. Focal object recognition **653** may indicate segmented, classified objects from a computer vision system or spatially indexed object database, together with gaze or gesture input to identify which object a viewing user is focusing on, or which is most salient to this user. Significant events detection **655** may encompass what is happening around a user in the scene **600**; public and private events may be detected or inferred based on the current activity, location, or context of a user. In particular embodiments, the context engine **510** may detect context changes, and trigger a series of events in response to relevant changes in downstream components, which may be registered to particular events in order to effect particular actions. As examples and not by way of limitation, such context changes may be people entering or exiting the scene, detection of a new object, determining that a person or object has been recorded for a threshold amount of time, movements to or from another movement type, starting or ending a particular activity, a user arriving or leaving a location, and detecting when a user has been in a current location for longer than a threshold amount of time.

With reference to FIG. 6A, when monitoring scene **600**, context engine **510** may detect people **602-605**, and determine the respective identity of each (e.g., Bob **602**, Alice **603**, Carol **604**, and Sarah **605**). These detected people may be entered into the “Social Presence” category **651** of chart **650**. Context engine **510** may further detect various activities performed by users in the scene **600** (e.g., walking, standing, eating, or speaking); these detected activities may be entered into the “User Activity Class” **652** of chart **650**. Context engine **510** may further detect objects **610-617** (e.g., cat **610**, coat **611**, water bottle **612**, table **613**, hat **614**, glasses **615**, and candy bowl **616**) and enter detected objects into the “Focal Object Recognition” category **653** of chart **650**. A location of the monitored scene **600**, such as an address, building, or room of the scene, may be determined and entered into the “User Location” category of chart **650**. A type of event, such as party occurring at the determined location, may further be detected, and entered into the “Significant Events Detection” category of chart **650**.

This information gathered by the context engine **510** may enable the assistant system **140** to execute various user- or object-specific functions. As an example and not by way of limitation, a user may make a request such as “where is Bob?”, and a camera of a client device at the viewed scene may locate and focus on Bob, as the context engine **510** has already identified which person is Bob. Similarly, a query of “who is speaking?” may result in this camera focusing on a speaker (e.g., Alice) while displaying the speaker’s name, and a request to take a picture of the speaker may activate the camera to take the requested photo. As another example and not by way of limitation, a request such as “follow the cat” may trigger the camera to locate and track the cat, as the context engine **510** has already identified various objects, including the cat, in scene **600**.

In particular embodiments, to answer more complex user questions and requests dealing with inter-object relationships, the assistant system **140** may access, from a scene-understanding engine of the assistant system **140**, relationship data associated with the video call. This relationship data may indicate relationships between various entities within the scene of the video call. In particular embodiments, the assistant system **140** may determine that the request references a particular type of relationship data. As an example and not by way of limitation, the request may

include a relationship word such as “holding” or “left of”, each of which indicates a different type of relationship among entities in the scene. In response to determining that the request references this particular type of relationship data, the assistant system **140** may activate scene understanding engine **520**. In particular embodiments, upon being activated, the scene understanding engine **520** may analyze the video call to generate relationship data of the particular type of relationship data referenced in the request. The scene understanding engine **520** may generate this relationship data in real time in response to being activated, and, after the relationship data has been generated, the assistant system **140** may deactivate the scene understanding engine **520**.

FIG. 6C illustrates an example knowledge graph **660** of the scene **600** generated by scene understanding engine **520**. While context engine **510** may always be on, the scene understanding engine **520** may be awakened as needed. The scene understanding engine **520** may receive data (such as the data of chart **650**) tracked by the context engine **510**, and determine relationships among the various detected entities, including both people and objects. As an example and not by way of limitation, the scene understanding engine **520** may determine that Sarah is holding the candy bowl (that the relationship between entities “Sarah” and “candy bowl” that have been identified by the context engine **510** is “holding”), that Sarah is looking at the cat, that the cat is carrying chicken, that Alice is to the left of Bob, etc. The scene understanding engine **520** may generate knowledge graph **660** of entity relationships; this knowledge graph **660** may be generated on and concern only the scene **600** of the video call. Because determining such semantic information may be computationally expensive, the scene understanding engine **520** may be awakened in response to a request that includes a relationship word (such as “holding” or “left of”), rather than remaining always on (or in ambient mode) like the context engine **510**. However, even in embodiments in which the scene understanding engine **520** awakens only upon request, the scene understanding engine **520** may be able to generate the information needed for a response relatively quickly using the specific information from the context engine **510** (for example, with respect to the question “what is Alice holding?”, the context engine **510** has already identified which person is Alice, and that the object is a water bottle). This relationship information output by the scene understanding engine **520** may enable the first user to ask otherwise ambiguous questions such as “what is Alice holding?” or “who is the guy wearing the blue hat?”, as well as commands such as “focus on the person next to Alice”. Further, the assistant system **140** may perform context carryover in order to answer a chain of such ambiguous questions. For example, a viewing user may ask “who is behind Bob?”, and the assistant system **140** may answer “Sarah is behind Bob”. Subsequently, the viewing user may ask “what is she looking at?”. The assistant system **140** may recognize that this question refers to the previously identified user “Sarah”, and respond with “Sarah is looking at a cat”.

In particular embodiments, the assistant system **140** may send, to one or more of the plurality of client systems **130**, a response to the request while maintaining the video call between the plurality of client systems **130**. In particular embodiments, which client systems(s) **130** the response is sent to may be based on the intent of the request. As an example and not by way of limitation, if the first user’s request was for information about the scene **600** of the video call (e.g., “who is speaking?”), or if the request was to virtually modify a characteristic (e.g., an appearance) of a

second user, the response may be sent to the first client system 130 of the first user. As another example and not by way of limitation, if the first user's request was to share content with one or more second users, the identified content may be sent to respective client systems 130 of those one or more second users. Although this disclosure describes sending responses to client systems 130 in a particular manner, this disclosure contemplates sending responses to client systems 130 in any suitable manner.

FIG. 7A illustrates an example shifted scene 700 viewed after a user command concerning an entity of the initial scene 600 on the first client system 130 of the first user 710. As an example and not by way of limitation, a command by user 710 to "follow the cat" may cause scene 600 to shift to the left as a smart camera on a client device of a user in the scene 600 tracks the cat to a different location, thus resulting in an updated scene 700. In particular embodiments, output of the context and scene understanding engines 510 and 520 may be dynamically updated as the viewed scene 600 changes. FIG. 7B illustrates an example updated chart 750 of information of the shifted scene 700 generated by the context engine 510. As can be seen in chart 750, new entities now visible in the shifted scene 700 such as Dave and Eve have been added to "Social Presence" category 651, while objects fireplace, cat bed, and wine glass have been added to "Focal Object Recognition" category 653. FIG. 7C illustrates an example updated knowledge graph 760 of the shifted scene 700 generated by the scene understanding engine 520. Updated knowledge graph 760 of the shifted scene 700 may be generated by the scene understanding engine 520 in response to a command containing a wake-word concerning a relationship between entities detected by the context engine 510, such as "what is the cat carrying?". Although the example knowledge graphs 660 and 760 illustrate particular words as the relationship information expressed by the edges between nodes representing entities and objects within a scene as determined by a context engine 510, each relationship edge may be represented by various synonyms or related words. For example, the relationship information "holding" may also map to requests involving "carrying", "lifting", or "having"; similarly, "right of" and "left of" may be considered types of a "next to" relationship.

FIG. 8 illustrates an example updated scene 800 viewed after a user command concerning an entity of a previous scene, such as scene 600, on the first client system 130 of the first user 810. As an example and not by way of limitation, effects and animations may be incorporated into the video call via the assistant system 140. User 810 may request that the assistant system 140 modify the appearance of a user (for example, upon noticing the cape-like quality of Carol's coat, user 810 may request that the assistant system 140 add special effects such as vampire fangs or wings to user Carol), or alter the voice of the user (for example, using the voice of another person in the chat). As another example and not by way of limitation, user 810 may request the assistant system 140 to add effects, such as a mask, onto another person, or to perform image processing to change the appearance of the other person (for example, making them appear to be younger or older than they are). The assistant system 140 may further incorporate AR/VR functionalities to allow the user 810 to request to, for example, give a user viewed in the scene 800 a virtual hug.

In particular embodiments, this architecture, combined with relevant social-networking information, may enable the assistant system 140 to determine the people and objects involved in the video call, thus enabling multiple social and sharing functions using these identifications. As an example

and not by way of limitation, user 810 may query the assistant system 140 as to where to buy the coat 611 that Carol is wearing, and the assistant system 140 may respond with a link to the appropriate product on an online shopping site.

FIG. 9 illustrates an example video call in which content relevant to the video call is viewed on the client system 130 of a user 910. As an example and not by way of limitation, during a video call 900, the assistant system 140 may provide content relevant to the video call or to a given user during a video call between users 910 and 920. For instance, the assistant system 140 may retrieve social-networking information of user 920 during the video call, and provide this information to a chatting user 910 to refer to during a conversation in the video call. In particular embodiments, the assistant system 140 may provide a content sidebar 930 for display, such as by causing the content sidebar to pop up or slide into the frame of the video call. Content sidebar 930 may present various content items, such as photos 931 and 932; these photos may be relevant to the two users 910 and 920 or to a topic or property of the call. As an example and not by way of limitation, photos 931 and 932 may be mutually tagged photos of the users. Content sidebar 930 may also present content items such as video 933; this video 933 may, similar to photos 931 and 932, involve users 910 and/or 920, or it may be a video of a subject that concerns a common interest of both users 910 and 920. Content sidebar 930 may also display social-networking content items, such as a post 934 that may be authored by one of the users 910 and 920 or a by mutual contact, or that may simply have one or more of the users 910 and 920 tagged in it. Content sidebar 930 may further present social-networking information 935, such as a list of mutual hobbies. Such content items 931-935 may increase the sense of social connection of the chatting users 910 and 920, and may further guide their conversation, for example, when it appears that they may be running out of topics to discuss.

Another technical challenge may be that, when conducting a video call on the client device 130, the user of that device may wish to preserve access to the functions of the device and access to a smart assistant system, which may go dormant during the video call. A solution presented by embodiments disclosed herein to address this challenge may thus involve the persistent assistant system 140 that, rather than going dormant during a video call, remains active but on standby, and is thus accessible to the first user to be invoked during the video call to execute various commands on the client device 130 of the first user. As another example, providing the persistent, always-on assistant system 140 may enable the first user to continue to use their client device 130 and assistant system 140 normally, even while conducting the video call.

In particular embodiments, the assistant system 140 may perform various social functions relating to users 910 and 920 during the video call. As an example and not by way of limitation, user 910 may query the assistant system 140 as to what a topic of conversation should be for a video call with another user 920, and, subject to privacy settings, the assistant system 140 may consult the social-networking information of that user 920 to identify common interests or a recent post created or shared by user 920. As another example and not by way of limitation, a request to share a "[song/video/picture] with them" may result in the relevant media being shared with the other user(s) 920 in the video call. In particular embodiments, the presence of the persistent, always-on assistant system 140 may enable user 910 to continue to use their client device 130 and assistant system

140 normally, even while conducting a video call. As an example and not by way of limitation, the assistant system 140 may be queried as to what a particular person just said. As another example and not by way of limitation, the assistant system 140 may take notes during the video call, or even summarize portions of or the entire call. As yet another example and not by way of limitation, with the persistent assistant system 140 to execute requests, user device functions such as timers, weather alerts, alarms, news, and other utilities that may ordinarily be dormant during a video call may remain accessible during that video call.

In particular embodiments, the persistent assistant system 140 and its resulting single communication domain may further enable various communication functionalities. As an example and not by way of limitation, as discussed above, a smart microphone integrated with the assistant system 140 may allow a user to mute himself while speaking to the assistant system 140, so that other users in the video call cannot hear him. As another example and not by way of limitation, the assistant system 140 may translate between speakers (e.g., if another person in the video call is speaking in Chinese, this speech may be translated in real time into English, and provided to the user either in audio or as captions). As yet another example and not by way of limitation, accessibility functions may also be enabled, such as interpreting a gaze or gesture of a user as input, rather than an explicit audio command.

FIG. 10 illustrates an example method 1000 for generating a response to a user request to a persistent assistant system 140 made during a video call. The method may begin at step 1010, where the persistent assistant system 140 may establish a video call between the client systems 130 of multiple users. At step 1020, the persistent assistant system 140 may receive a request from a first user of a first client system 130 to be performed by the persistent assistant system 140 during the video call. At step 1030, the persistent assistant system 140 may determine user identifiers of second users referenced in the request, as well as determining the intent of the request. For example, a user viewing a scene may command the persistent assistant system 140 to take a photo of a second user within the scene with a camera, follow a certain user, or show where a second user is located within the scene; the assistant system 140 may then determine the user identifier of this user in order to perform the requested action. At step 1040, the assistant system 140 may be instructed to execute the request based on the user identifiers and intent. At step 1050, the assistant system 140 may send a response to the request to one or more of the second users while maintaining the video call. For example, if the request was to take a photo of a second user, the smart camera may zoom in on the identified second user and take the requested photo, and then display that photo to one or more participants on the video call during the actual video call.

Particular embodiments may repeat one or more steps of the method of FIG. 10, where appropriate. Although this disclosure describes and illustrates particular steps of the method of FIG. 10 as occurring in a particular order, this disclosure contemplates any suitable steps of the method of FIG. 10 occurring in any suitable order. Moreover, although this disclosure describes and illustrates an example method for generating a response to a user request made to a persistent assistant system during a video call including the particular steps of the method of FIG. 10, this disclosure contemplates any suitable method for generating a response to a user request made to a persistent assistant system during a video call, including any suitable steps, which may include

all, some, or none of the steps of the method of FIG. 10, where appropriate. Furthermore, although this disclosure describes and illustrates particular components, devices, or systems carrying out particular steps of the method of FIG. 10, this disclosure contemplates any suitable combination of any suitable components, devices, or systems carrying out any suitable steps of the method of FIG. 10.

Social Graphs

FIG. 11 illustrates an example social graph 1100. In particular embodiments, the social-networking system 160 may store one or more social graphs 1100 in one or more data stores. In particular embodiments, the social graph 1100 may include multiple nodes—which may include multiple user nodes 1102 or multiple concept nodes 1104—and multiple edges 1106 connecting the nodes. Each node may be associated with a unique entity (i.e., user or concept), each of which may have a unique identifier (ID), such as a unique number or username. The example social graph 1100 illustrated in FIG. 11 is shown, for didactic purposes, in a two-dimensional visual map representation. In particular embodiments, a social-networking system 160, a client system 130, an assistant system 140, or a third-party system 170 may access the social graph 1100 and related social-graph information for suitable applications. The nodes and edges of the social graph 1100 may be stored as data objects, for example, in a data store (such as a social-graph database). Such a data store may include one or more searchable or queryable indexes of nodes or edges of the social graph 1100.

In particular embodiments, a user node 1102 may correspond to a user of the social-networking system 160 or the assistant system 140. As an example and not by way of limitation, a user may be an individual (human user), an entity (e.g., an enterprise, business, or third-party application), or a group (e.g., of individuals or entities) that interacts or communicates with or over the social-networking system 160 or the assistant system 140. In particular embodiments, when a user registers for an account with the social-networking system 160, the social-networking system 160 may create a user node 1102 corresponding to the user, and store the user node 1102 in one or more data stores. Users and user nodes 1102 described herein may, where appropriate, refer to registered users and user nodes 1102 associated with registered users. In addition or as an alternative, users and user nodes 1102 described herein may, where appropriate, refer to users that have not registered with the social-networking system 160. In particular embodiments, a user node 1102 may be associated with information provided by a user or information gathered by various systems, including the social-networking system 160. As an example and not by way of limitation, a user may provide his or her name, profile picture, contact information, birth date, sex, marital status, family status, employment, education background, preferences, interests, or other demographic information. In particular embodiments, a user node 1102 may be associated with one or more data objects corresponding to information associated with a user. In particular embodiments, a user node 1102 may correspond to one or more web interfaces.

In particular embodiments, a concept node 1104 may correspond to a concept. As an example and not by way of limitation, a concept may correspond to a place (such as, for example, a movie theater, restaurant, landmark, or city); a website (such as, for example, a website associated with the social-networking system 160 or a third-party website associated with a web-application server); an entity (such as, for example, a person, business, group, sports team, or celebrity); a resource (such as, for example, an audio file, video

file, digital photo, text file, structured document, or application) which may be located within the social-networking system **160** or on an external server, such as a web-application server; real or intellectual property (such as, for example, a sculpture, painting, movie, game, song, idea, photograph, or written work); a game; an activity; an idea or theory; another suitable concept; or two or more such concepts. A concept node **1104** may be associated with information of a concept provided by a user or information gathered by various systems, including the social-networking system **160** and the assistant system **140**. As an example and not by way of limitation, information of a concept may include a name or a title; one or more images (e.g., an image of the cover page of a book); a location (e.g., an address or a geographical location); a website (which may be associated with a URL); contact information (e.g., a phone number or an email address); other suitable concept information; or any suitable combination of such information. In particular embodiments, a concept node **1104** may be associated with one or more data objects corresponding to information associated with concept node **1104**. In particular embodiments, a concept node **1104** may correspond to one or more web interfaces.

In particular embodiments, a node in the social graph **1100** may represent or be represented by a web interface (which may be referred to as a “profile interface”). Profile interfaces may be hosted by or accessible to the social-networking system **160** or the assistant system **140**. Profile interfaces may also be hosted on third-party websites associated with a third-party system **170**. As an example and not by way of limitation, a profile interface corresponding to a particular external web interface may be the particular external web interface and the profile interface may correspond to a particular concept node **1104**. Profile interfaces may be viewable by all or a selected subset of other users. As an example and not by way of limitation, a user node **1102** may have a corresponding user-profile interface in which the corresponding user may add content, make declarations, or otherwise express himself or herself. As another example and not by way of limitation, a concept node **1104** may have a corresponding concept-profile interface in which one or more users may add content, make declarations, or express themselves, particularly in relation to the concept corresponding to concept node **1104**.

In particular embodiments, a concept node **1104** may represent a third-party web interface or resource hosted by a third-party system **170**. The third-party web interface or resource may include, among other elements, content, a selectable or other icon, or other inter-actable object representing an action or activity. As an example and not by way of limitation, a third-party web interface may include a selectable icon such as “like,” “check-in,” “eat,” “recommend,” or another suitable action or activity. A user viewing the third-party web interface may perform an action by selecting one of the icons (e.g., “check-in”), causing a client system **130** to send to the social-networking system **160** a message indicating the user’s action. In response to the message, the social-networking system **160** may create an edge (e.g., a check-in-type edge) between a user node **1102** corresponding to the user and a concept node **1104** corresponding to the third-party web interface or resource and store edge **1106** in one or more data stores.

In particular embodiments, a pair of nodes in the social graph **1100** may be connected to each other by one or more edges **1106**. An edge **1106** connecting a pair of nodes may represent a relationship between the pair of nodes. In particular embodiments, an edge **1106** may include or represent

one or more data objects or attributes corresponding to the relationship between a pair of nodes. As an example and not by way of limitation, a first user may indicate that a second user is a “friend” of the first user. In response to this indication, the social-networking system **160** may send a “friend request” to the second user. If the second user confirms the “friend request,” the social-networking system **160** may create an edge **1106** connecting the first user’s user node **1102** to the second user’s user node **1102** in the social graph **1100** and store edge **1106** as social-graph information in one or more of data stores **164**. In the example of FIG. **11**, the social graph **1100** includes an edge **1106** indicating a friend relation between user nodes **1102** of user “A” and user “B” and an edge indicating a friend relation between user nodes **1102** of user “C” and user “B.” Although this disclosure describes or illustrates particular edges **1106** with particular attributes connecting particular user nodes **1102**, this disclosure contemplates any suitable edges **1106** with any suitable attributes connecting user nodes **1102**. As an example and not by way of limitation, an edge **1106** may represent a friendship, family relationship, business or employment relationship, fan relationship (including, e.g., liking, etc.), follower relationship, visitor relationship (including, e.g., accessing, viewing, checking-in, sharing, etc.), subscriber relationship, superior/subordinate relationship, reciprocal relationship, non-reciprocal relationship, another suitable type of relationship, or two or more such relationships. Moreover, although this disclosure generally describes nodes as being connected, this disclosure also describes users or concepts as being connected. Herein, references to users or concepts being connected may, where appropriate, refer to the nodes corresponding to those users or concepts being connected in the social graph **1100** by one or more edges **1106**. The degree of separation between two objects represented by two nodes, respectively, is a count of edges in a shortest path connecting the two nodes in the social graph **1100**. As an example and not by way of limitation, in the social graph **1100**, the user node **1102** of user “C” is connected to the user node **1102** of user “A” via multiple paths including, for example, a first path directly passing through the user node **1102** of user “B,” a second path passing through the concept node **1104** of company “Alme” and the user node **1102** of user “D,” and a third path passing through the user nodes **1102** and concept nodes **1104** representing school “Stateford,” user “G,” company “Alme,” and user “D.” User “C” and user “A” have a degree of separation of two because the shortest path connecting their corresponding nodes (i.e., the first path) includes two edges **1106**.

In particular embodiments, an edge **1106** between a user node **1102** and a concept node **1104** may represent a particular action or activity performed by a user associated with user node **1102** toward a concept associated with a concept node **1104**. As an example and not by way of limitation, as illustrated in FIG. **11**, a user may “like,” “attended,” “played,” “listened,” “cooked,” “worked at,” or “read” a concept, each of which may correspond to an edge type or subtype. A concept-profile interface corresponding to a concept node **1104** may include, for example, a selectable “check in” icon (such as, for example, a clickable “check in” icon) or a selectable “add to favorites” icon. Similarly, after a user clicks these icons, the social-networking system **160** may create a “favorite” edge or a “check in” edge in response to a user’s action corresponding to a respective action. As another example and not by way of limitation, a user (user “C”) may listen to a particular song (“Imagine”) using a particular application (a third-party online music

application). In this case, the social-networking system 160 may create a “listened” edge 1106 and a “used” edge (as illustrated in FIG. 11) between user nodes 1102 corresponding to the user and concept nodes 1104 corresponding to the song and application to indicate that the user listened to the song and used the application. Moreover, the social-networking system 160 may create a “played” edge 1106 (as illustrated in FIG. 11) between concept nodes 1104 corresponding to the song and the application to indicate that the particular song was played by the particular application. In this case, “played” edge 1106 corresponds to an action performed by an external application (the third-party online music application) on an external audio file (the song “Imagine”). Although this disclosure describes particular edges 1106 with particular attributes connecting user nodes 1102 and concept nodes 1104, this disclosure contemplates any suitable edges 1106 with any suitable attributes connecting user nodes 1102 and concept nodes 1104. Moreover, although this disclosure describes edges between a user node 1102 and a concept node 1104 representing a single relationship, this disclosure contemplates edges between a user node 1102 and a concept node 1104 representing one or more relationships. As an example and not by way of limitation, an edge 1106 may represent both that a user likes and has used at a particular concept. Alternatively, another edge 1106 may represent each type of relationship (or multiples of a single relationship) between a user node 1102 and a concept node 1104 (as illustrated in FIG. 11 between user node 1102 for user “E” and concept node 1104 for “online music application”).

In particular embodiments, the social-networking system 160 may create an edge 1106 between a user node 1102 and a concept node 1104 in the social graph 1100. As an example and not by way of limitation, a user viewing a concept-profile interface (such as, for example, by using a web browser or a special-purpose application hosted by the user’s client system 130) may indicate that he or she likes the concept represented by the concept node 1104 by clicking or selecting a “Like” icon, which may cause the user’s client system 130 to send to the social-networking system 160 a message indicating the user’s liking of the concept associated with the concept-profile interface. In response to the message, the social-networking system 160 may create an edge 1106 between user node 1102 associated with the user and concept node 1104, as illustrated by “like” edge 1106 between the user and concept node 1104. In particular embodiments, the social-networking system 160 may store an edge 1106 in one or more data stores. In particular embodiments, an edge 1106 may be automatically formed by the social-networking system 160 in response to a particular user action. As an example and not by way of limitation, if a first user uploads a picture, reads a book, watches a movie, or listens to a song, an edge 1106 may be formed between user node 1102 corresponding to the first user and concept nodes 1104 corresponding to those concepts. Although this disclosure describes forming particular edges 1106 in particular manners, this disclosure contemplates forming any suitable edges 1106 in any suitable manner.

Vector Spaces and Embeddings

FIG. 12 illustrates an example view of a vector space 1200. In particular embodiments, an object or an n-gram may be represented in a d-dimensional vector space, where d denotes any suitable number of dimensions. Although the vector space 1200 is illustrated as a three-dimensional space, this is for illustrative purposes only, as the vector space 1200 may be of any suitable dimension. In particular embodi-

ments, an n-gram may be represented in the vector space 1200 as a vector referred to as a term embedding. Each vector may comprise coordinates corresponding to a particular point in the vector space 1200 (i.e., the terminal point of the vector). As an example and not by way of limitation, vectors 1210, 1220, and 1230 may be represented as points in the vector space 1200, as illustrated in FIG. 12. An n-gram may be mapped to a respective vector representation. As an example and not by way of limitation, n-grams t_1 and t_2 may be mapped to vectors \vec{v}_1 and \vec{v}_2 in the vector space 1200, respectively, by applying a function $\vec{\pi}$ defined by a dictionary, such that $\vec{v}_1 = \vec{\pi}(t_1)$ and $\vec{v}_2 = \vec{\pi}(t_2)$. As another example and not by way of limitation, a dictionary trained to map text to a vector representation may be utilized, or such a dictionary may be itself generated via training. As another example and not by way of limitation, a word-embeddings model may be used to map an n-gram to a vector representation in the vector space 1200. In particular embodiments, an n-gram may be mapped to a vector representation in the vector space 1200 by using a machine learning model (e.g., a neural network). The machine learning model may have been trained using a sequence of training data (e.g., a corpus of objects each comprising n-grams).

In particular embodiments, an object may be represented in the vector space 1200 as a vector referred to as a feature vector or an object embedding. As an example and not by way of limitation, objects e_1 and e_2 may be mapped to vectors \vec{v}_1 and \vec{v}_2 in the vector space 1200, respectively, by applying a function $\vec{\pi}$, such that $\vec{v}_1 = \vec{\pi}(e_1)$ and $\vec{v}_2 = \vec{\pi}(e_2)$. In particular embodiments, an object may be mapped to a vector based on one or more properties, attributes, or features of the object, relationships of the object with other objects, or any other suitable information associated with the object. As an example and not by way of limitation, a function $\vec{\pi}$ may map objects to vectors by feature extraction, which may start from an initial set of measured data and build derived values (e.g., features). As an example and not by way of limitation, an object comprising a video or an image may be mapped to a vector by using an algorithm to detect or isolate various desired portions or shapes of the object. Features used to calculate the vector may be based on information obtained from edge detection, corner detection, blob detection, ridge detection, scale-invariant feature transformation, edge direction, changing intensity, autocorrelation, motion detection, optical flow, thresholding, blob extraction, template matching, Hough transformation (e.g., lines, circles, ellipses, arbitrary shapes), or any other suitable information. As another example and not by way of limitation, an object comprising audio data may be mapped to a vector based on features such as a spectral slope, a tonality coefficient, an audio spectrum centroid, an audio spectrum envelope, a Mel-frequency cepstrum, or any other suitable information. In particular embodiments, when an object has data that is either too large to be efficiently processed or comprises redundant data, a function $\vec{\pi}$ may map the object to a vector using a transformed reduced set of features (e.g., feature selection). In particular embodiments, a function $\vec{\pi}$ may map an object e to a vector $\vec{\pi}(e)$ based on one or more n-grams associated with object e . Although this disclosure describes representing an n-gram or an object in a vector space in a particular manner, this disclosure contemplates representing an n-gram or an object in a vector space in any suitable manner.

In particular embodiments, the social-networking system 160 may calculate a similarity metric of vectors in vector space 1200. A similarity metric may be a cosine similarity, a Minkowski distance, a Mahalanobis distance, a Jaccard

similarity coefficient, or any suitable similarity metric. As an example and not by way of limitation, a similarity metric \vec{v}_1 of and \vec{v}_2 may be a cosine similarity

$$\frac{\vec{v}_1 \cdot \vec{v}_2}{\|\vec{v}_1\| \|\vec{v}_2\|}$$

As another example and not by way of limitation, a similarity metric of \vec{v}_1 and \vec{v}_2 may be a Euclidean distance $\|\vec{v}_1 - \vec{v}_2\|$. A similarity metric of two vectors may represent how similar the two objects or n-grams corresponding to the two vectors, respectively, are to one another, as measured by the distance between the two vectors in the vector space **1200**. As an example and not by way of limitation, vector **1210** and vector **1220** may correspond to objects that are more similar to one another than the objects corresponding to vector **1210** and vector **1230**, based on the distance between the respective vectors. Although this disclosure describes calculating a similarity metric between vectors in a particular manner, this disclosure contemplates calculating a similarity metric between vectors in any suitable manner.

More information on vector spaces, embeddings, feature vectors, and similarity metrics may be found in U.S. patent application Ser. No. 14/949,436, filed 23 Nov. 2015, U.S. patent application Ser. No. 15/286,315, filed 5 Oct. 2016, and U.S. patent application Ser. No. 15/365,789, filed 30 Nov. 2016, each of which is incorporated by reference. Artificial Neural Networks

FIG. **13** illustrates an example artificial neural network (“ANN”) **1300**. In particular embodiments, an ANN may refer to a computational model comprising one or more nodes. Example ANN **1300** may comprise an input layer **1310**, hidden layers **1320**, **1330**, **1340**, and an output layer **1350**. Each layer of the ANN **1300** may comprise one or more nodes, such as a node **1305** or a node **1315**. In particular embodiments, each node of an ANN may be connected to another node of the ANN. As an example and not by way of limitation, each node of the input layer **1310** may be connected to one of more nodes of the hidden layer **1320**. In particular embodiments, one or more nodes may be a bias node (e.g., a node in a layer that is not connected to and does not receive input from any node in a previous layer). In particular embodiments, each node in each layer may be connected to one or more nodes of a previous or subsequent layer. Although FIG. **13** depicts a particular ANN with a particular number of layers, a particular number of nodes, and particular connections between nodes, this disclosure contemplates any suitable ANN with any suitable number of layers, any suitable number of nodes, and any suitable connections between nodes. As an example and not by way of limitation, although FIG. **13** depicts a connection between each node of the input layer **1310** and each node of the hidden layer **1320**, one or more nodes of the input layer **1310** may not be connected to one or more nodes of the hidden layer **1320**.

In particular embodiments, an ANN may be a feedforward ANN (e.g., an ANN with no cycles or loops where communication between nodes flows in one direction beginning with the input layer and proceeding to successive layers). As an example and not by way of limitation, the input to each node of the hidden layer **1320** may comprise the output of one or more nodes of the input layer **1310**. As another example and not by way of limitation, the input to each node of the output layer **1350** may comprise the output of one or

more nodes of the hidden layer **1340**. In particular embodiments, an ANN may be a deep neural network (e.g., a neural network comprising at least two hidden layers). In particular embodiments, an ANN may be a deep residual network. A deep residual network may be a feedforward ANN comprising hidden layers organized into residual blocks. The input into each residual block after the first residual block may be a function of the output of the previous residual block and the input of the previous residual block. As an example and not by way of limitation, the input into residual block N may be $F(x)+x$, where $F(x)$ may be the output of residual block N-1, x may be the input into residual block N-1. Although this disclosure describes a particular ANN, this disclosure contemplates any suitable ANN.

In particular embodiments, an activation function may correspond to each node of an ANN. An activation function of a node may define the output of a node for a given input. In particular embodiments, an input to a node may comprise a set of inputs. As an example and not by way of limitation, an activation function may be an identity function, a binary step function, a logistic function, or any other suitable function. As another example and not by way of limitation, an activation function for a node k may be the sigmoid function

$$F_k(s_k) = \frac{1}{1 + e^{-s_k}}$$

the hyperbolic tangent function

$$F_k(s_k) = \frac{e^{s_k} - e^{-s_k}}{e^{s_k} + e^{-s_k}}$$

the rectifier $F_k(s_k) = \max(0, s_k)$, or any other suitable function $F_k(s_k)$, where s_k may be the effective input to node k. In particular embodiments, the input of an activation function corresponding to a node may be weighted. Each node may generate output using a corresponding activation function based on weighted inputs. In particular embodiments, each connection between nodes may be associated with a weight. As an example and not by way of limitation, a connection **1325** between the node **1305** and the node **1315** may have a weighting coefficient of 0.4, which may indicate that 0.4 multiplied by the output of the node **1305** is used as an input to the node **1315**. As another example and not by way of limitation, the output y_k of node k may be $y_k = F_k(s_k)$, where F_k may be the activation function corresponding to node k, $s_k = \sum_j (w_{jk} x_j)$ may be the effective input to node k, x_j may be the output of a node j connected to node k, and w_{jk} may be the weighting coefficient between node j and node k. In particular embodiments, the input to nodes of the input layer may be based on a vector representing an object. Although this disclosure describes particular inputs to and outputs of nodes, this disclosure contemplates any suitable inputs to and outputs of nodes. Moreover, although this disclosure may describe particular connections and weights between nodes, this disclosure contemplates any suitable connections and weights between nodes.

In particular embodiments, an ANN may be trained using training data. As an example and not by way of limitation, training data may comprise inputs to the ANN **1300** and an expected output. As another example and not by way of limitation, training data may comprise vectors each representing a training object and an expected label for each training object. In particular embodiments, training an ANN

may comprise modifying the weights associated with the connections between nodes of the ANN by optimizing an objective function. As an example and not by way of limitation, a training method may be used (e.g., the conjugate gradient method, the gradient descent method, the stochastic gradient descent) to backpropagate the sum-of-squares error measured as a distances between each vector representing a training object (e.g., using a cost function that minimizes the sum-of-squares error). In particular embodiments, an ANN may be trained using a dropout technique. As an example and not by way of limitation, one or more nodes may be temporarily omitted (e.g., receive no input and generate no output) while training. For each training object, one or more nodes of the ANN may have some probability of being omitted. The nodes that are omitted for a particular training object may be different than the nodes omitted for other training objects (e.g., the nodes may be temporarily omitted on an object-by-object basis). Although this disclosure describes training an ANN in a particular manner, this disclosure contemplates training an ANN in any suitable manner.

Privacy

In particular embodiments, one or more objects (e.g., content or other types of objects) of a computing system may be associated with one or more privacy settings. The one or more objects may be stored on or otherwise associated with any suitable computing system or application, such as, for example, a social-networking system **160**, a client system **130**, an assistant system **140**, a third-party system **170**, a social-networking application, an assistant application, a messaging application, a photo-sharing application, or any other suitable computing system or application. Although the examples discussed herein are in the context of an online social network, these privacy settings may be applied to any other suitable computing system. Privacy settings (or “access settings”) for an object may be stored in any suitable manner, such as, for example, in association with the object, in an index on an authorization server, in another suitable manner, or any suitable combination thereof. A privacy setting for an object may specify how the object (or particular information associated with the object) can be accessed, stored, or otherwise used (e.g., viewed, shared, modified, copied, executed, surfaced, or identified) within the online social network. When privacy settings for an object allow a particular user or other entity to access that object, the object may be described as being “visible” with respect to that user or other entity. As an example and not by way of limitation, a user of the online social network may specify privacy settings for a user-profile page that identify a set of users that may access work-experience information on the user-profile page, thus excluding other users from accessing that information.

In particular embodiments, privacy settings for an object may specify a “blocked list” of users or other entities that should not be allowed to access certain information associated with the object. In particular embodiments, the blocked list may include third-party entities. The blocked list may specify one or more users or entities for which an object is not visible. As an example and not by way of limitation, a user may specify a set of users who may not access photo albums associated with the user, thus excluding those users from accessing the photo albums (while also possibly allowing certain users not within the specified set of users to access the photo albums). In particular embodiments, privacy settings may be associated with particular social-graph elements. Privacy settings of a social-graph element, such as a node or an edge, may specify how the social-graph

element, information associated with the social-graph element, or objects associated with the social-graph element can be accessed using the online social network. As an example and not by way of limitation, a particular concept node **1104** corresponding to a particular photo may have a privacy setting specifying that the photo may be accessed only by users tagged in the photo and friends of the users tagged in the photo. In particular embodiments, privacy settings may allow users to opt in to or opt out of having their content, information, or actions stored/logged by the social-networking system **160** or assistant system **140** or shared with other systems (e.g., a third-party system **170**). Although this disclosure describes using particular privacy settings in a particular manner, this disclosure contemplates using any suitable privacy settings in any suitable manner.

In particular embodiments, privacy settings may be based on one or more nodes or edges of a social graph **1100**. A privacy setting may be specified for one or more edges **1106** or edge-types of the social graph **1100**, or with respect to one or more nodes **1102**, **1104** or node-types of the social graph **1100**. The privacy settings applied to a particular edge **1106** connecting two nodes may control whether the relationship between the two entities corresponding to the nodes is visible to other users of the online social network. Similarly, the privacy settings applied to a particular node may control whether the user or concept corresponding to the node is visible to other users of the online social network. As an example and not by way of limitation, a first user may share an object to the social-networking system **160**. The object may be associated with a concept node **1104** connected to a user node **1102** of the first user by an edge **1106**. The first user may specify privacy settings that apply to a particular edge **1106** connecting to the concept node **1104** of the object, or may specify privacy settings that apply to all edges **1106** connecting to the concept node **1104**. As another example and not by way of limitation, the first user may share a set of objects of a particular object-type (e.g., a set of images). The first user may specify privacy settings with respect to all objects associated with the first user of that particular object-type as having a particular privacy setting (e.g., specifying that all images posted by the first user are visible only to friends of the first user and/or users tagged in the images).

In particular embodiments, the social-networking system **160** may present a “privacy wizard” (e.g., within a webpage, a module, one or more dialog boxes, or any other suitable interface) to the first user to assist the first user in specifying one or more privacy settings. The privacy wizard may display instructions, suitable privacy-related information, current privacy settings, one or more input fields for accepting one or more inputs from the first user specifying a change or confirmation of privacy settings, or any suitable combination thereof. In particular embodiments, the social-networking system **160** may offer a “dashboard” functionality to the first user that may display, to the first user, current privacy settings of the first user. The dashboard functionality may be displayed to the first user at any appropriate time (e.g., following an input from the first user summoning the dashboard functionality, following the occurrence of a particular event or trigger action). The dashboard functionality may allow the first user to modify one or more of the first user’s current privacy settings at any time, in any suitable manner (e.g., redirecting the first user to the privacy wizard).

Privacy settings associated with an object may specify any suitable granularity of permitted access or denial of access. As an example and not by way of limitation, access or denial of access may be specified for particular users (e.g.,

only me, my roommates, my boss), users within a particular degree-of-separation (e.g., friends, friends-of-friends), user groups (e.g., the gaming club, my family), user networks (e.g., employees of particular employers, students or alumni of particular university), all users (“public”), no users (“private”), users of third-party systems **170**, particular applications (e.g., third-party applications, external websites), other suitable entities, or any suitable combination thereof. Although this disclosure describes particular granularities of permitted access or denial of access, this disclosure contemplates any suitable granularities of permitted access or denial of access.

In particular embodiments, one or more servers **162** may be authorization/privacy servers for enforcing privacy settings. In response to a request from a user (or other entity) for a particular object stored in a data store **164**, the social-networking system **160** may send a request to the data store **164** for the object. The request may identify the user associated with the request and the object may be sent only to the user (or a client system **130** of the user) if the authorization server determines that the user is authorized to access the object based on the privacy settings associated with the object. If the requesting user is not authorized to access the object, the authorization server may prevent the requested object from being retrieved from the data store **164** or may prevent the requested object from being sent to the user. In the search-query context, an object may be provided as a search result only if the querying user is authorized to access the object, e.g., if the privacy settings for the object allow it to be surfaced to, discovered by, or otherwise visible to the querying user. In particular embodiments, an object may represent content that is visible to a user through a newsfeed of the user. As an example and not by way of limitation, one or more objects may be visible to a user’s “Trending” page. In particular embodiments, an object may correspond to a particular user. The object may be content associated with the particular user, or may be the particular user’s account or information stored on the social-networking system **160**, or other computing system. As an example and not by way of limitation, a first user may view one or more second users of an online social network through a “People You May Know” function of the online social network, or by viewing a list of friends of the first user. As an example and not by way of limitation, a first user may specify that they do not wish to see objects associated with a particular second user in their newsfeed or friends list. If the privacy settings for the object do not allow it to be surfaced to, discovered by, or visible to the user, the object may be excluded from the search results. Although this disclosure describes enforcing privacy settings in a particular manner, this disclosure contemplates enforcing privacy settings in any suitable manner.

In particular embodiments, different objects of the same type associated with a user may have different privacy settings. Different types of objects associated with a user may have different types of privacy settings. As an example and not by way of limitation, a first user may specify that the first user’s status updates are public, but any images shared by the first user are visible only to the first user’s friends on the online social network. As another example and not by way of limitation, a user may specify different privacy settings for different types of entities, such as individual users, friends-of-friends, followers, user groups, or corporate entities. As another example and not by way of limitation, a first user may specify a group of users that may view videos posted by the first user, while keeping the videos from being visible to the first user’s employer. In particular

embodiments, different privacy settings may be provided for different user groups or user demographics. As an example and not by way of limitation, a first user may specify that other users who attend the same university as the first user may view the first user’s pictures, but that other users who are family members of the first user may not view those same pictures.

In particular embodiments, the social-networking system **160** may provide one or more default privacy settings for each object of a particular object-type. A privacy setting for an object that is set to a default may be changed by a user associated with that object. As an example and not by way of limitation, all images posted by a first user may have a default privacy setting of being visible only to friends of the first user and, for a particular image, the first user may change the privacy setting for the image to be visible to friends and friends-of-friends.

In particular embodiments, privacy settings may allow a first user to specify (e.g., by opting out, by not opting in) whether the social-networking system **160** or assistant system **140** may receive, collect, log, or store particular objects or information associated with the user for any purpose. In particular embodiments, privacy settings may allow the first user to specify whether particular applications or processes may access, store, or use particular objects or information associated with the user. The privacy settings may allow the first user to opt in or opt out of having objects or information accessed, stored, or used by specific applications or processes. The social-networking system **160** or assistant system **140** may access such information in order to provide a particular function or service to the first user, without the social-networking system **160** or assistant system **140** having access to that information for any other purposes. Before accessing, storing, or using such objects or information, the social-networking system **160** or assistant system **140** may prompt the user to provide privacy settings specifying which applications or processes, if any, may access, store, or use the object or information prior to allowing any such action. As an example and not by way of limitation, a first user may transmit a message to a second user via an application related to the online social network (e.g., a messaging app), and may specify privacy settings that such messages should not be stored by the social-networking system **160** or assistant system **140**.

In particular embodiments, a user may specify whether particular types of objects or information associated with the first user may be accessed, stored, or used by the social-networking system **160** or assistant system **140**. As an example and not by way of limitation, the first user may specify that images sent by the first user through the social-networking system **160** or assistant system **140** may not be stored by the social-networking system **160** or assistant system **140**. As another example and not by way of limitation, a first user may specify that messages sent from the first user to a particular second user may not be stored by the social-networking system **160** or assistant system **140**. As yet another example and not by way of limitation, a first user may specify that all objects sent via a particular application may be saved by the social-networking system **160** or assistant system **140**.

In particular embodiments, privacy settings may allow a first user to specify whether particular objects or information associated with the first user may be accessed from particular client systems **130** or third-party systems **170**. The privacy settings may allow the first user to opt in or opt out of having objects or information accessed from a particular device (e.g., the phone book on a user’s smart phone), from

a particular application (e.g., a messaging app), or from a particular system (e.g., an email server). The social-networking system **160** or assistant system **140** may provide default privacy settings with respect to each device, system, or application, and/or the first user may be prompted to specify a particular privacy setting for each context. As an example and not by way of limitation, the first user may utilize a location-services feature of the social-networking system **160** or assistant system **140** to provide recommendations for restaurants or other places in proximity to the user. The first user's default privacy settings may specify that the social-networking system **160** or assistant system **140** may use location information provided from a client device **130** of the first user to provide the location-based services, but that the social-networking system **160** or assistant system **140** may not store the location information of the first user or provide it to any third-party system **170**. The first user may then update the privacy settings to allow location information to be used by a third-party image-sharing application in order to geo-tag photos.

In particular embodiments, privacy settings may allow a user to specify one or more geographic locations from which objects can be accessed. Access or denial of access to the objects may depend on the geographic location of a user who is attempting to access the objects. As an example and not by way of limitation, a user may share an object and specify that only users in the same city may access or view the object. As another example and not by way of limitation, a first user may share an object and specify that the object is visible to second users only while the first user is in a particular location. If the first user leaves the particular location, the object may no longer be visible to the second users. As another example and not by way of limitation, a first user may specify that an object is visible only to second users within a threshold distance from the first user. If the first user subsequently changes location, the original second users with access to the object may lose access, while a new group of second users may gain access as they come within the threshold distance of the first user.

In particular embodiments, the social-networking system **160** or assistant system **140** may have functionalities that may use, as inputs, personal or biometric information of a user for user-authentication or experience-personalization purposes. A user may opt to make use of these functionalities to enhance their experience on the online social network. As an example and not by way of limitation, a user may provide personal or biometric information to the social-networking system **160** or assistant system **140**. The user's privacy settings may specify that such information may be used only for particular processes, such as authentication, and further specify that such information may not be shared with any third-party system **170** or used for other processes or applications associated with the social-networking system **160** or assistant system **140**. As another example and not by way of limitation, the social-networking system **160** may provide a functionality for a user to provide voice-print recordings to the online social network. As an example and not by way of limitation, if a user wishes to utilize this function of the online social network, the user may provide a voice recording of his or her own voice to provide a status update on the online social network. The recording of the voice-input may be compared to a voice print of the user to determine what words were spoken by the user. The user's privacy setting may specify that such voice recording may be used only for voice-input purposes (e.g., to authenticate the user, to send voice messages, to improve voice recognition in order to use voice-operated features of the online social network), and

further specify that such voice recording may not be shared with any third-party system **170** or used by other processes or applications associated with the social-networking system **160**. As another example and not by way of limitation, the social-networking system **160** may provide a functionality for a user to provide a reference image (e.g., a facial profile, a retinal scan) to the online social network. The online social network may compare the reference image against a later-received image input (e.g., to authenticate the user, to tag the user in photos). The user's privacy setting may specify that such image may be used only for a limited purpose (e.g., authentication, tagging the user in photos), and further specify that such image may not be shared with any third-party system **170** or used by other processes or applications associated with the social-networking system **160**.

Systems and Methods

FIG. **14** illustrates an example computer system **1400**. In particular embodiments, one or more computer systems **1400** perform one or more steps of one or more methods described or illustrated herein. In particular embodiments, one or more computer systems **1400** provide functionality described or illustrated herein. In particular embodiments, software running on one or more computer systems **1400** performs one or more steps of one or more methods described or illustrated herein or provides functionality described or illustrated herein. Particular embodiments include one or more portions of one or more computer systems **1400**. Herein, reference to a computer system may encompass a computing device, and vice versa, where appropriate. Moreover, reference to a computer system may encompass one or more computer systems, where appropriate.

This disclosure contemplates any suitable number of computer systems **1400**. This disclosure contemplates computer system **1400** taking any suitable physical form. As an example and not by way of limitation, computer system **1400** may be an embedded computer system, a system-on-chip (SOC), a single-board computer system (SBC) (such as, for example, a computer-on-module (COM) or system-on-module (SOM)), a desktop computer system, a laptop or notebook computer system, an interactive kiosk, a mainframe, a mesh of computer systems, a mobile telephone, a personal digital assistant (PDA), a server, a tablet computer system, or a combination of two or more of these. Where appropriate, computer system **1400** may include one or more computer systems **1400**; be unitary or distributed; span multiple locations; span multiple machines; span multiple data centers; or reside in a cloud, which may include one or more cloud components in one or more networks. Where appropriate, one or more computer systems **1400** may perform without substantial spatial or temporal limitation one or more steps of one or more methods described or illustrated herein. As an example and not by way of limitation, one or more computer systems **1400** may perform in real time or in batch mode one or more steps of one or more methods described or illustrated herein. One or more computer systems **1400** may perform at different times or at different locations one or more steps of one or more methods described or illustrated herein, where appropriate.

In particular embodiments, computer system **1400** includes a processor **1402**, memory **1404**, storage **1406**, an input/output (I/O) interface **1408**, a communication interface **1410**, and a bus **1412**. Although this disclosure describes and illustrates a particular computer system having a particular number of particular components in a particular arrangement, this disclosure contemplates any suitable computer

system having any suitable number of any suitable components in any suitable arrangement.

In particular embodiments, processor **1402** includes hardware for executing instructions, such as those making up a computer program. As an example and not by way of limitation, to execute instructions, processor **1402** may retrieve (or fetch) the instructions from an internal register, an internal cache, memory **1404**, or storage **1406**; decode and execute them; and then write one or more results to an internal register, an internal cache, memory **1404**, or storage **1406**. In particular embodiments, processor **1402** may include one or more internal caches for data, instructions, or addresses. This disclosure contemplates processor **1402** including any suitable number of any suitable internal caches, where appropriate. As an example and not by way of limitation, processor **1402** may include one or more instruction caches, one or more data caches, and one or more translation lookaside buffers (TLBs). Instructions in the instruction caches may be copies of instructions in memory **1404** or storage **1406**, and the instruction caches may speed up retrieval of those instructions by processor **1402**. Data in the data caches may be copies of data in memory **1404** or storage **1406** for instructions executing at processor **1402** to operate on; the results of previous instructions executed at processor **1402** for access by subsequent instructions executing at processor **1402** or for writing to memory **1404** or storage **1406**; or other suitable data. The data caches may speed up read or write operations by processor **1402**. The TLBs may speed up virtual-address translation for processor **1402**. In particular embodiments, processor **1402** may include one or more internal registers for data, instructions, or addresses. This disclosure contemplates processor **1402** including any suitable number of any suitable internal registers, where appropriate. Where appropriate, processor **1402** may include one or more arithmetic logic units (ALUs); be a multi-core processor; or include one or more processors **1402**. Although this disclosure describes and illustrates a particular processor, this disclosure contemplates any suitable processor.

In particular embodiments, memory **1404** includes main memory for storing instructions for processor **1402** to execute or data for processor **1402** to operate on. As an example and not by way of limitation, computer system **1400** may load instructions from storage **1406** or another source (such as, for example, another computer system **1400**) to memory **1404**. Processor **1402** may then load the instructions from memory **1404** to an internal register or internal cache. To execute the instructions, processor **1402** may retrieve the instructions from the internal register or internal cache and decode them. During or after execution of the instructions, processor **1402** may write one or more results (which may be intermediate or final results) to the internal register or internal cache. Processor **1402** may then write one or more of those results to memory **1404**. In particular embodiments, processor **1402** executes only instructions in one or more internal registers or internal caches or in memory **1404** (as opposed to storage **1406** or elsewhere) and operates only on data in one or more internal registers or internal caches or in memory **1404** (as opposed to storage **1406** or elsewhere). One or more memory buses (which may each include an address bus and a data bus) may couple processor **1402** to memory **1404**. Bus **1412** may include one or more memory buses, as described below. In particular embodiments, one or more memory management units (MMUs) reside between processor **1402** and memory **1404** and facilitate accesses to memory **1404** requested by processor **1402**. In particular embodiments, memory **1404**

includes random access memory (RAM). This RAM may be volatile memory, where appropriate. Where appropriate, this RAM may be dynamic RAM (DRAM) or static RAM (SRAM). Moreover, where appropriate, this RAM may be single-ported or multi-ported RAM. This disclosure contemplates any suitable RAM. Memory **1404** may include one or more memories **1404**, where appropriate. Although this disclosure describes and illustrates particular memory, this disclosure contemplates any suitable memory.

In particular embodiments, storage **1406** includes mass storage for data or instructions. As an example and not by way of limitation, storage **1406** may include a hard disk drive (HDD), a floppy disk drive, flash memory, an optical disc, a magneto-optical disc, magnetic tape, or a Universal Serial Bus (USB) drive or a combination of two or more of these. Storage **1406** may include removable or non-removable (or fixed) media, where appropriate. Storage **1406** may be internal or external to computer system **1400**, where appropriate. In particular embodiments, storage **1406** is non-volatile, solid-state memory. In particular embodiments, storage **1406** includes read-only memory (ROM). Where appropriate, this ROM may be mask-programmed ROM, programmable ROM (PROM), erasable PROM (EPROM), electrically erasable PROM (EEPROM), electrically alterable ROM (EAROM), or flash memory or a combination of two or more of these. This disclosure contemplates mass storage **1406** taking any suitable physical form. Storage **1406** may include one or more storage control units facilitating communication between processor **1402** and storage **1406**, where appropriate. Where appropriate, storage **1406** may include one or more storages **1406**. Although this disclosure describes and illustrates particular storage, this disclosure contemplates any suitable storage.

In particular embodiments, I/O interface **1408** includes hardware, software, or both, providing one or more interfaces for communication between computer system **1400** and one or more I/O devices. Computer system **1400** may include one or more of these I/O devices, where appropriate. One or more of these I/O devices may enable communication between a person and computer system **1400**. As an example and not by way of limitation, an I/O device may include a keyboard, keypad, microphone, monitor, mouse, printer, scanner, speaker, still camera, stylus, tablet, touch screen, trackball, video camera, another suitable I/O device or a combination of two or more of these. An I/O device may include one or more sensors. This disclosure contemplates any suitable I/O devices and any suitable I/O interfaces **1408** for them. Where appropriate, I/O interface **1408** may include one or more device or software drivers enabling processor **1402** to drive one or more of these I/O devices. I/O interface **1408** may include one or more I/O interfaces **1408**, where appropriate. Although this disclosure describes and illustrates a particular I/O interface, this disclosure contemplates any suitable I/O interface.

In particular embodiments, communication interface **1410** includes hardware, software, or both providing one or more interfaces for communication (such as, for example, packet-based communication) between computer system **1400** and one or more other computer systems **1400** or one or more networks. As an example and not by way of limitation, communication interface **1410** may include a network interface controller (NIC) or network adapter for communicating with an Ethernet or other wire-based network or a wireless NIC (WNIC) or wireless adapter for communicating with a wireless network, such as a WI-FI network. This disclosure contemplates any suitable network and any suitable communication interface **1410** for it. As an

example and not by way of limitation, computer system **1400** may communicate with an ad hoc network, a personal area network (PAN), a local area network (LAN), a wide area network (WAN), a metropolitan area network (MAN), or one or more portions of the Internet or a combination of two or more of these. One or more portions of one or more of these networks may be wired or wireless. As an example, computer system **1400** may communicate with a wireless PAN (WPAN) (such as, for example, a BLUETOOTH WPAN), a WI-FI network, a WI-MAX network, a cellular telephone network (such as, for example, a Global System for Mobile Communications (GSM) network), or other suitable wireless network or a combination of two or more of these. Computer system **1400** may include any suitable communication interface **1410** for any of these networks, where appropriate. Communication interface **1410** may include one or more communication interfaces **1410**, where appropriate. Although this disclosure describes and illustrates a particular communication interface, this disclosure contemplates any suitable communication interface.

In particular embodiments, bus **1412** includes hardware, software, or both coupling components of computer system **1400** to each other. As an example and not by way of limitation, bus **1412** may include an Accelerated Graphics Port (AGP) or other graphics bus, an Enhanced Industry Standard Architecture (EISA) bus, a front-side bus (FSB), a HYPERTRANSPORT (HT) interconnect, an Industry Standard Architecture (ISA) bus, an INFINIBAND interconnect, a low-pin-count (LPC) bus, a memory bus, a Micro Channel Architecture (MCA) bus, a Peripheral Component Interconnect (PCI) bus, a PCI-Express (PCIe) bus, a serial advanced technology attachment (SATA) bus, a Video Electronics Standards Association local (VLB) bus, or another suitable bus or a combination of two or more of these. Bus **1412** may include one or more buses **1412**, where appropriate. Although this disclosure describes and illustrates a particular bus, this disclosure contemplates any suitable bus or interconnect.

Herein, a computer-readable non-transitory storage medium or media may include one or more semiconductor-based or other integrated circuits (ICs) (such as, for example, field-programmable gate arrays (FPGAs) or application-specific ICs (ASICs)), hard disk drives (HDDs), hybrid hard drives (HHDs), optical discs, optical disc drives (ODDs), magneto-optical discs, magneto-optical drives, floppy diskettes, floppy disk drives (FDDs), magnetic tapes, solid-state drives (SSDs), RAM-drives, SECURE DIGITAL cards or drives, any other suitable computer-readable non-transitory storage media, or any suitable combination of two or more of these, where appropriate. A computer-readable non-transitory storage medium may be volatile, non-volatile, or a combination of volatile and non-volatile, where appropriate.

Miscellaneous

Herein, “or” is inclusive and not exclusive, unless expressly indicated otherwise or indicated otherwise by context. Therefore, herein, “A or B” means “A, B, or both,” unless expressly indicated otherwise or indicated otherwise by context. Moreover, “and” is both joint and several, unless expressly indicated otherwise or indicated otherwise by context. Therefore, herein, “A and B” means “A and B, jointly or severally,” unless expressly indicated otherwise or indicated otherwise by context.

The scope of this disclosure encompasses all changes, substitutions, variations, alterations, and modifications to the example embodiments described or illustrated herein that a person having ordinary skill in the art would comprehend.

The scope of this disclosure is not limited to the example embodiments described or illustrated herein. Moreover, although this disclosure describes and illustrates respective embodiments herein as including particular components, elements, feature, functions, operations, or steps, any of these embodiments may include any combination or permutation of any of the components, elements, features, functions, operations, or steps described or illustrated anywhere herein that a person having ordinary skill in the art would comprehend. Furthermore, reference in the appended claims to an apparatus or system or a component of an apparatus or system being adapted to, arranged to, capable of, configured to, enabled to, operable to, or operative to perform a particular function encompasses that apparatus, system, component, whether or not it or that particular function is activated, turned on, or unlocked, as long as that apparatus, system, or component is so adapted, arranged, capable, configured, enabled, operable, or operative. Additionally, although this disclosure describes or illustrates particular embodiments as providing particular advantages, particular embodiments may provide none, some, or all of these advantages.

What is claimed is:

1. A method comprising, by one or more computing systems:
 - establishing a video call between a plurality of client systems, wherein access to an assistant system is persistently maintained during the video call;
 - receiving, from a first client system of the plurality of client systems, a request by a first user to be performed by the assistant system during the video call, wherein the request references one or more second users associated with the plurality of client systems;
 - determining an intent of the request;
 - analyzing, by a context engine of the assistant system, images of a scene of the video call to identify the one or more second users within the scene;
 - determining respective user identifiers of the one or more second users identified within the scene;
 - instructing the assistant system to execute the request based on the determined intent and user identifiers; and
 - sending, to one more of the plurality of client systems, a response to the request while maintaining the video call between the plurality of client systems.
2. The method of claim 1, wherein receiving the request comprises:
 - detecting a gaze of the first user, wherein the gaze is directed at one or more of the second users; and
 - inferring the request based on the detected gaze.
3. The method of claim 1, wherein the request comprises an instruction to focus a display of the first client system on one or more of the second users.
4. The method of claim 1, wherein the request comprises an instruction to repeat or summarize speech of one or more of the second users.
5. The method of claim 1, wherein the request comprises an instruction to share a content item with one or more of the second users.
6. The method of claim 1, further comprising:
 - retrieving user profile information of one or more of the second users in response to the determined intent and the one or more user identifiers; and
 - generating the response based on the retrieved user profile information.
7. The method of claim 6, wherein the user profile information comprises information of an interest or a recent activity of one or more of the second users.

59

8. The method of claim 1, wherein:
the determined intent is to modify an appearance of one
or more of the second users; and
the response comprises instructions to display, on the first
client system, a modified appearance of one or more of
the second users corresponding to the determined
intent.
9. The method of claim 1, wherein the request comprises
an instruction to perform a virtual activity with respect to
one or more of the second users.
10. The method of claim 1, wherein identifying informa-
tion of the first user is concealed during the video call.
11. The method of claim 1, further comprising:
receiving a wake-word that precedes the request; and
sending, in response to receiving the wake-word, instruc-
tions for muting the video call at the first client system.
12. The method of claim 1, further comprising:
accessing, from a context engine of the assistant system,
context data associated with the video call, wherein the
context data indicates properties of a scene of the video
call.
13. The method of claim 12, wherein the context data
comprises identifications of objects within the scene.
14. The method of claim 12, wherein the context engine
analyzes properties of the scene in real time during the video
call.
15. The method of claim 12, further comprising:
accessing, from a scene understanding engine of the
assistant system, relationship data associated with the
video call, wherein the relationship data indicates rela-
tionships between entities within the scene of the video
call.
16. The method of claim 15, further comprising:
determining the request references a particular type of
relationship data; and
activating the scene understanding engine in response to
determining the request references the particular type
of relationship data,
wherein the scene understanding engine analyzes the
video call to generate relationship data of the particular
type of relationship data referenced in the request.
17. The method of claim 16, wherein the scene under-
standing engine generates the relationship data in real time
in response to being activated.
18. The method of claim 16, further comprising:
deactivating the scene understanding engine after the
relationship data has been generated.

60

19. One or more computer-readable non-transitory stor-
age media embodying software that is operable when
executed to:
establish a video call between a plurality of client sys-
tems, wherein access to an assistant system is persis-
tently maintained during the video call;
receive, from a first client system of the plurality of client
systems, a request by a first user to be performed by the
assistant system during the video call, wherein the
request references one or more second users associated
with the plurality of client systems;
determine an intent of the request;
analyze, by a context engine of the assistant system,
images of a scene of the video call to identify the one
or more second users within the scene;
determine respective user identifiers of the one or more
second users identified within the scene;
instruct the assistant system to execute the request based
on the determined intent and user identifiers; and
send, to one more of the plurality of client systems, a
response to the request while maintaining the video call
between the plurality of client systems.
20. A system comprising one or more processors and one
or more computer-readable non-transitory storage media
coupled to one or more of the processors and comprising
instructions operable, when executed by one or more of the
processors, to cause the system to:
establish a video call between a plurality of client sys-
tems, wherein access to an assistant system is persis-
tently maintained during the video call;
receive, from a first client system of the plurality of client
systems, a request by a first user to be performed by the
assistant system during the video call, wherein the
request references one or more second users associated
with the plurality of client systems;
determine an intent of the request;
analyze, by a context engine of the assistant system,
images of a scene of the video call to identify the one
or more second users within the scene;
determine respective user identifiers of the one or more
second users identified within the scene;
instruct the assistant system to execute the request based
on the determined intent and user identifiers; and
send, to one more of the plurality of client systems, a
response to the request while maintaining the video call
between the plurality of client systems.

* * * * *